

DR

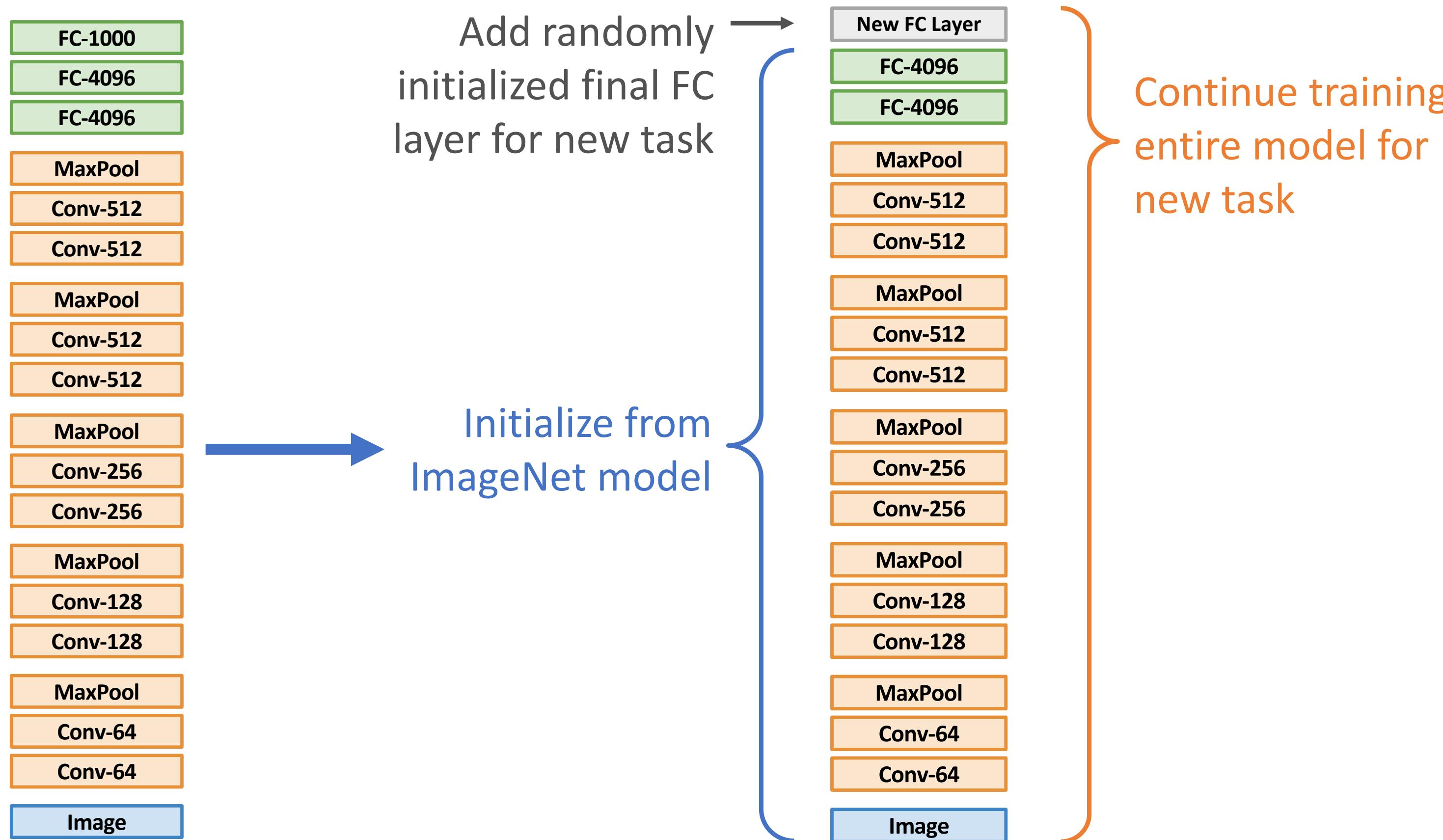
DeepRob

Lecture 13
Object Detectors and Segmentation
University of Michigan and University of Minnesota



Last time: Transfer Learning

1. Train on ImageNet



Last time: Localization Tasks

Classification



“Chocolate Pretzels”

No spatial extent



Semantic Segmentation



Chocolate Pretzels,
Shelf

No objects, just pixels

Object Detection



Flipz, Hershey's, Keese's

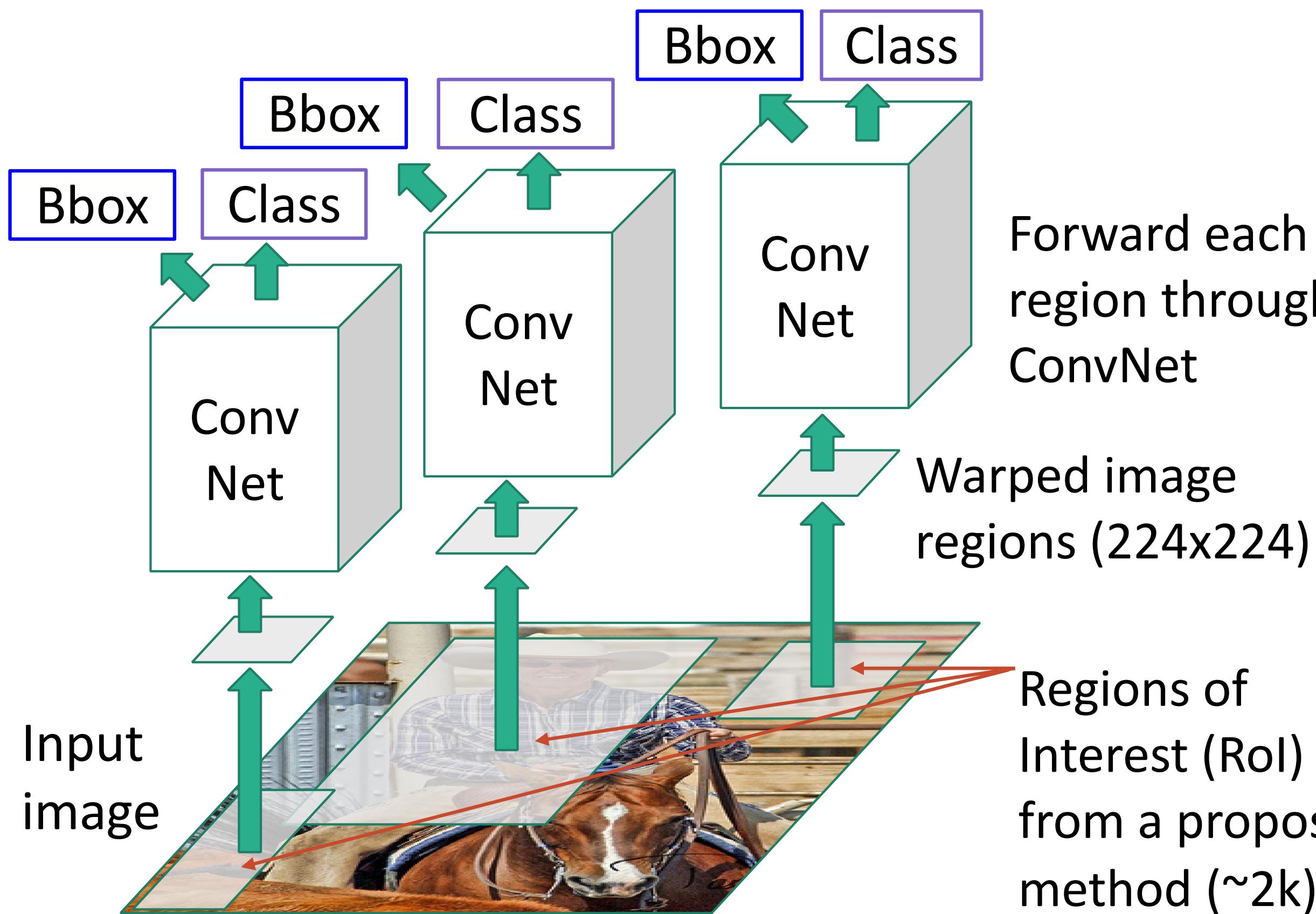
Multiple objects

Instance Segmentation



Last time: R-CNN

R-CNN: Region-Based CNN



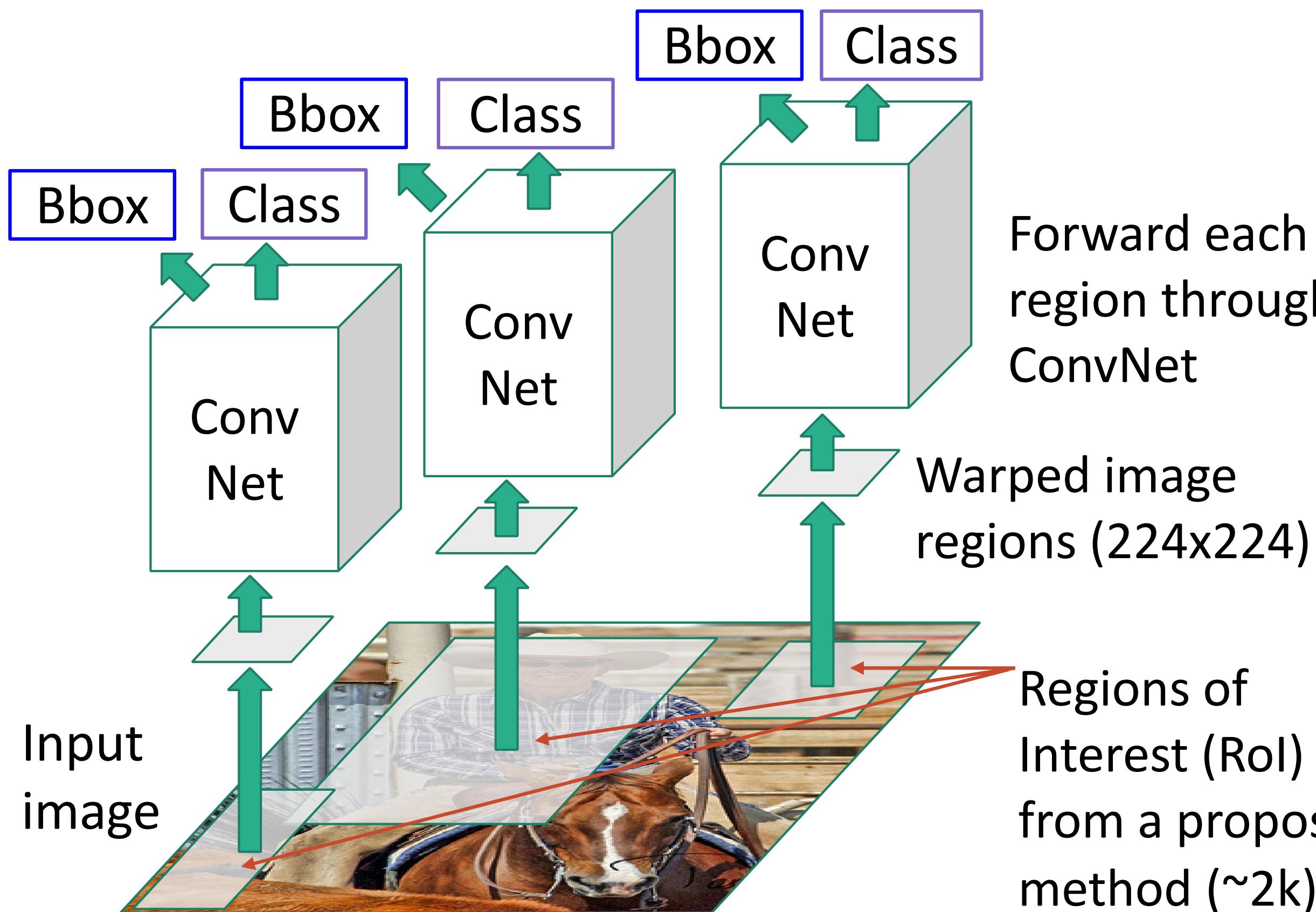
Classify each region

Bounding box regression:
Predict “transform” to correct the RoI: 4 numbers (t_x, t_y, t_h, t_w)

Problem: Very slow! Need to do 2000 forward passes through CNN per image

Last time: R-CNN

R-CNN: Region-Based CNN



Classify each region

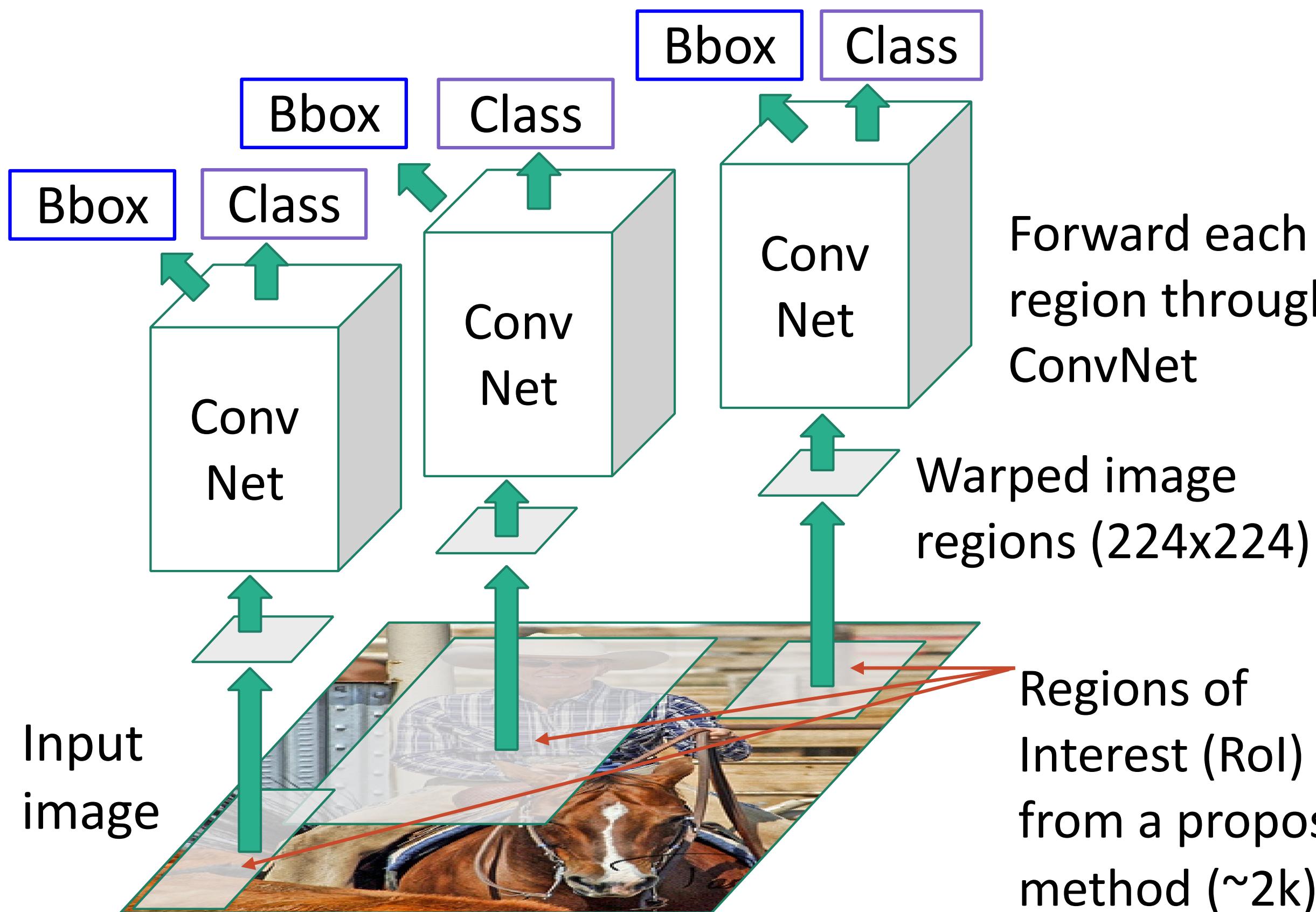
Bounding box regression:
Predict “transform” to correct the RoI: 4 numbers (t_x, t_y, t_h, t_w)

Problem: Very slow! Need to do 2000 forward passes through CNN per image

Idea: Overlapping proposals cause a lot of repeated work; same pixels processed many times. Can we avoid this?

Last time: R-CNN

R-CNN: Region-Based CNN



Classify each region

Bounding box regression:
Predict “transform” to correct the RoI: 4 numbers (t_x, t_y, t_h, t_w)

Problem: Very slow! Need to do 2000 forward passes through CNN per image

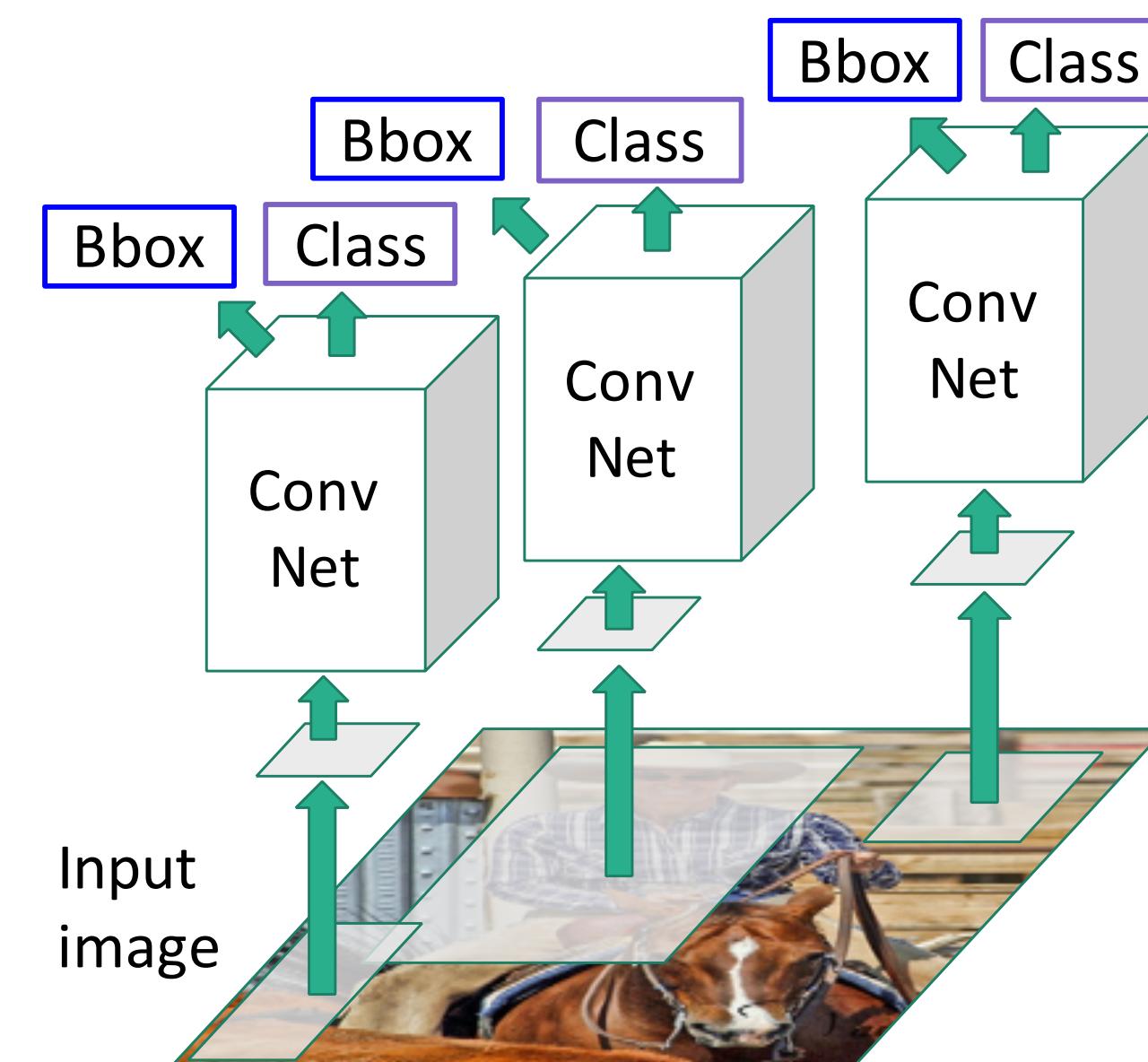
Idea: Overlapping proposals cause a lot of repeated work; same pixels processed many times. Can we avoid this?

Fast R-CNN



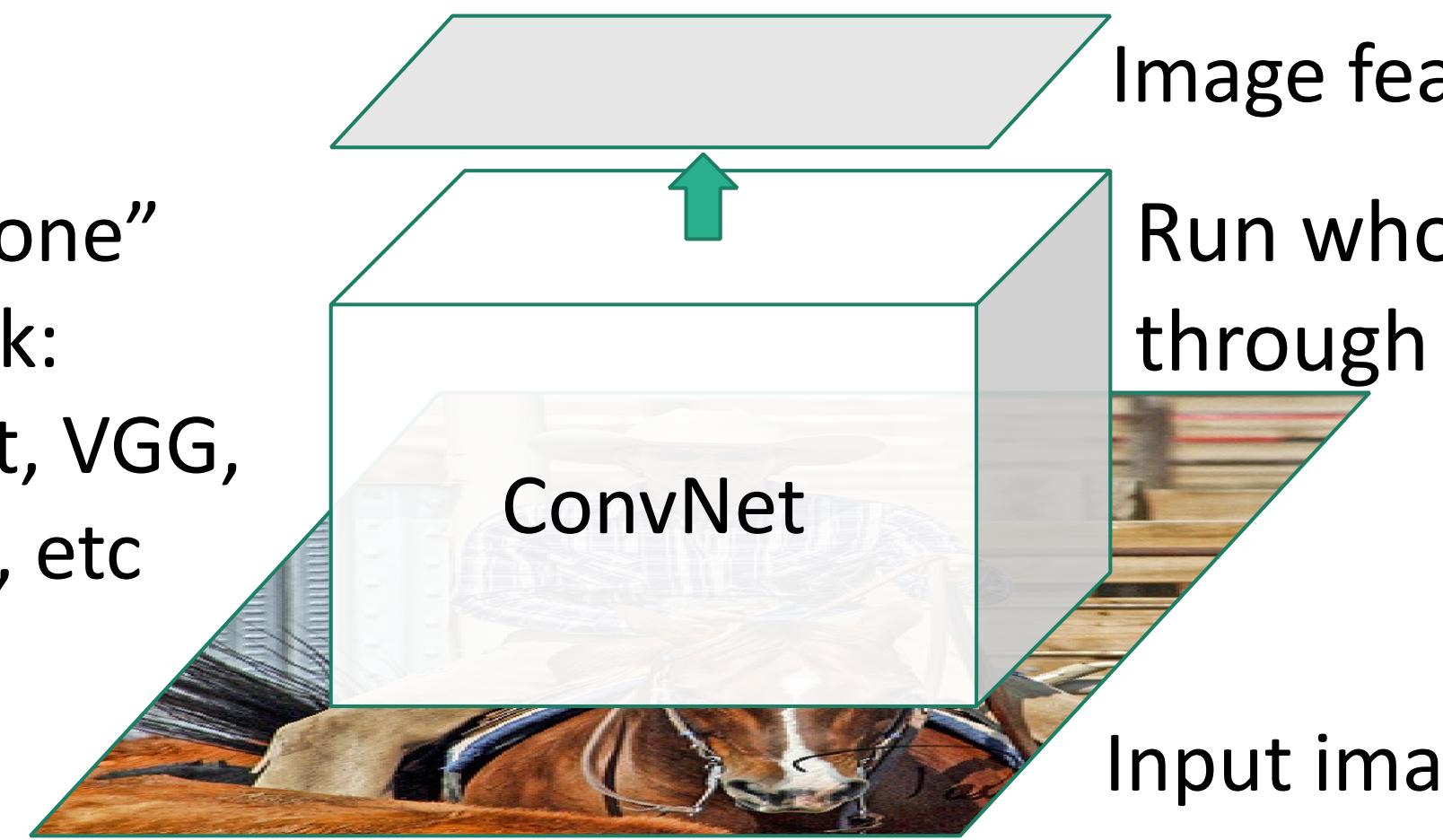
Input image

“Slow” R-CNN
Process each region
independently



Fast R-CNN

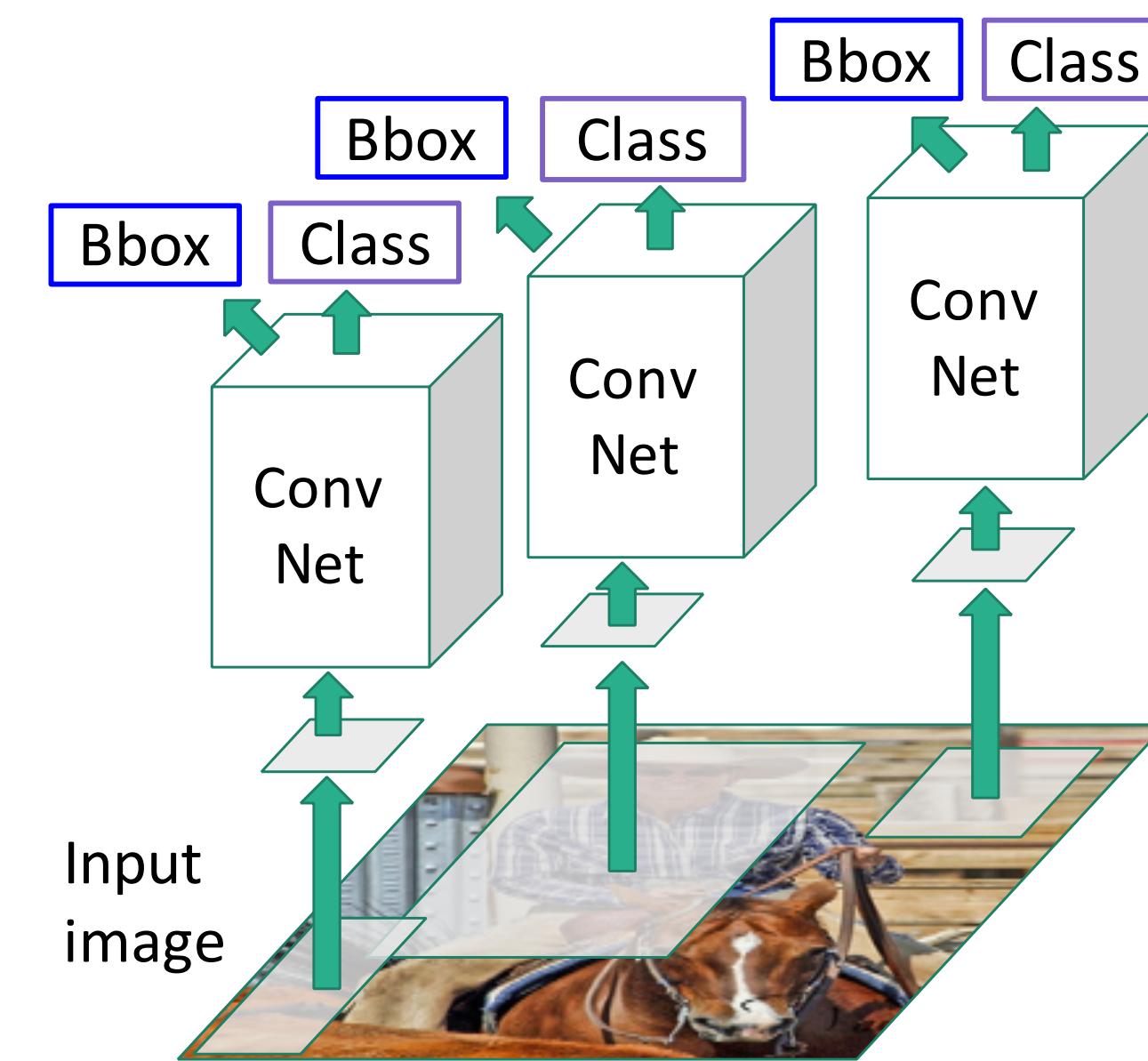
“Backbone” network:
AlexNet, VGG,
ResNet, etc



Run whole image
through ConvNet

Input image

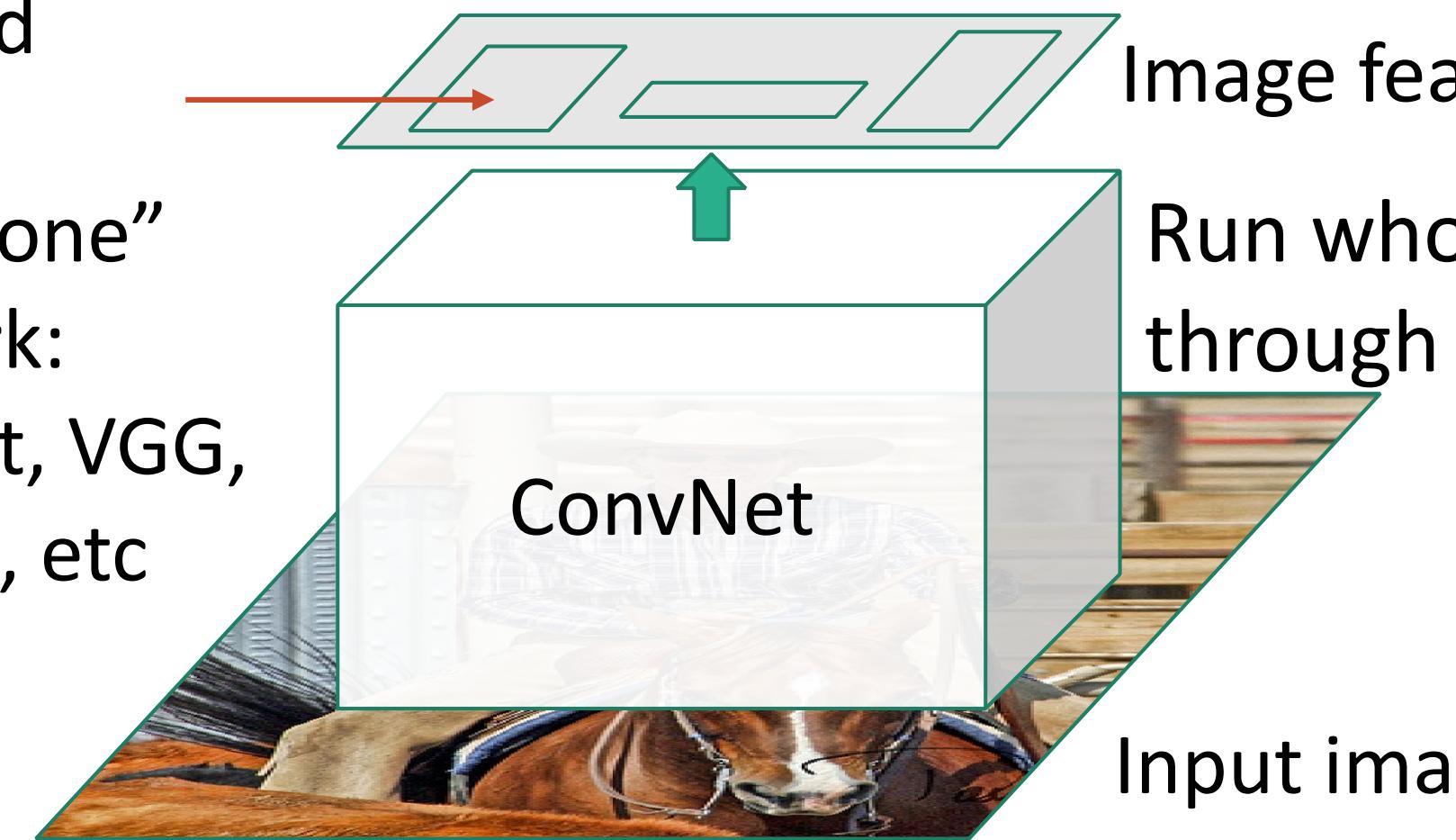
“Slow” R-CNN
Process each region
independently



Fast R-CNN

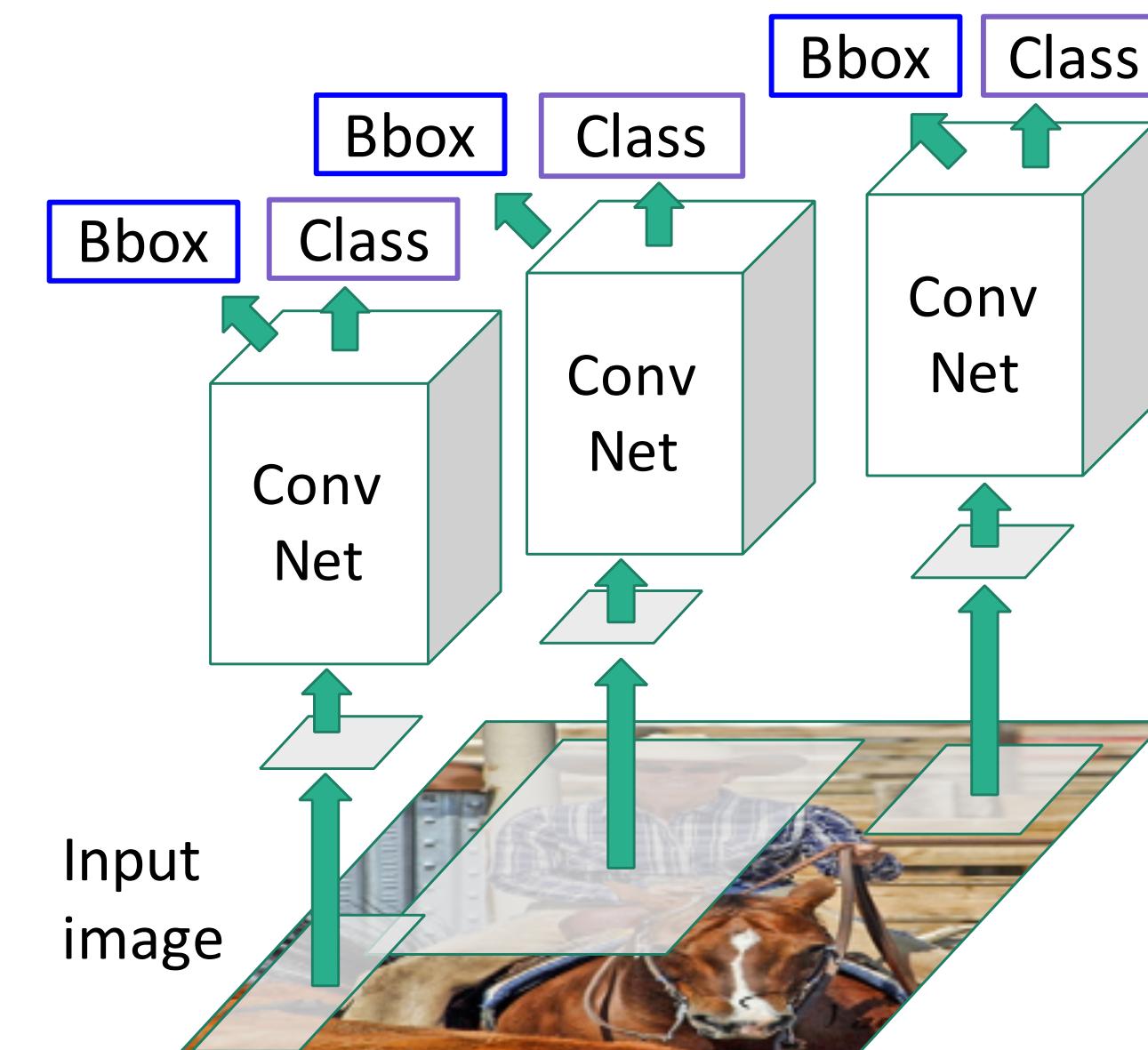
Regions of Interest (Rois)
from a proposal
method

“Backbone”
network:
AlexNet, VGG,
ResNet, etc



Run whole image
through ConvNet
Input image

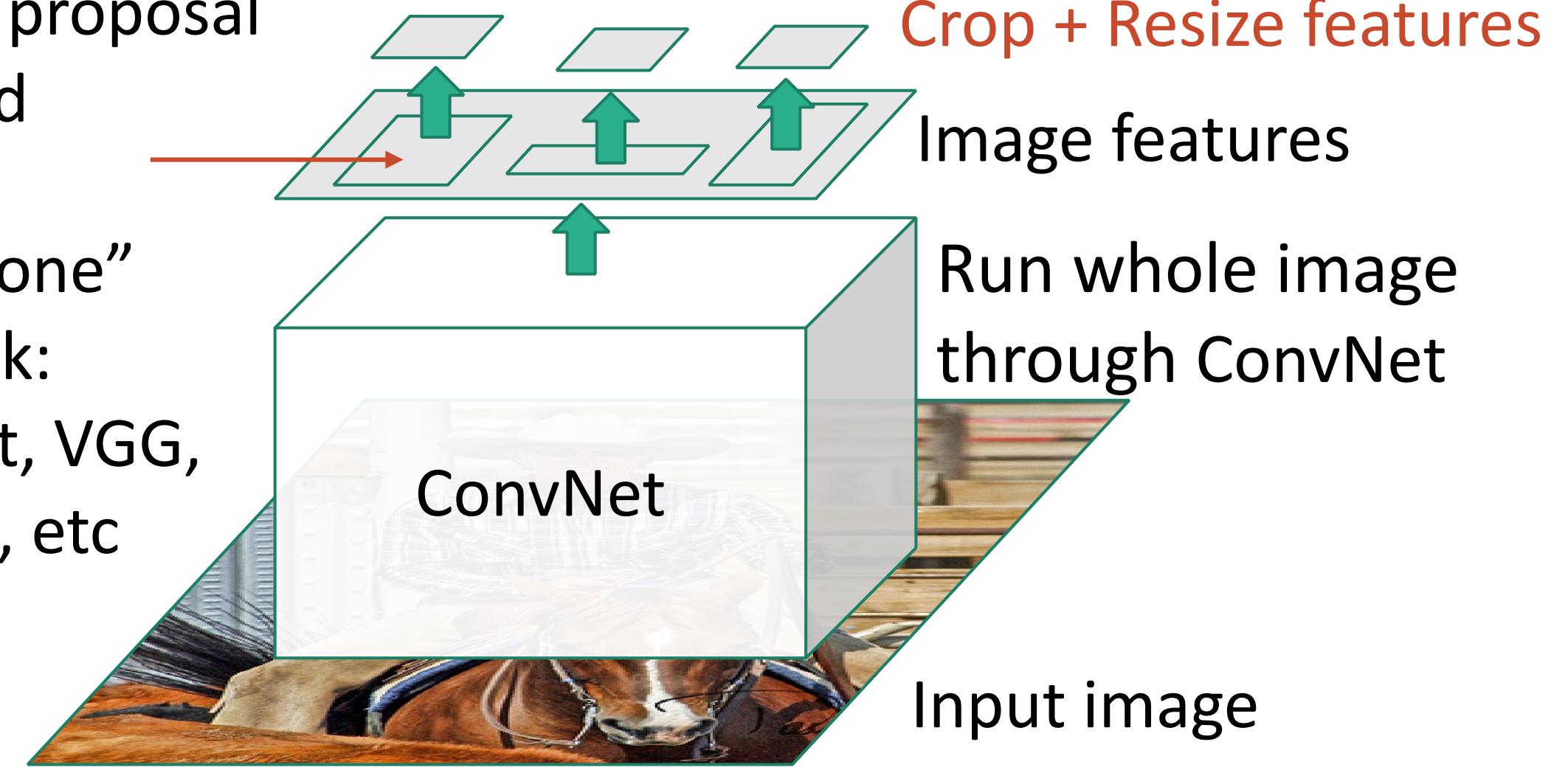
“Slow” R-CNN
Process each region
independently



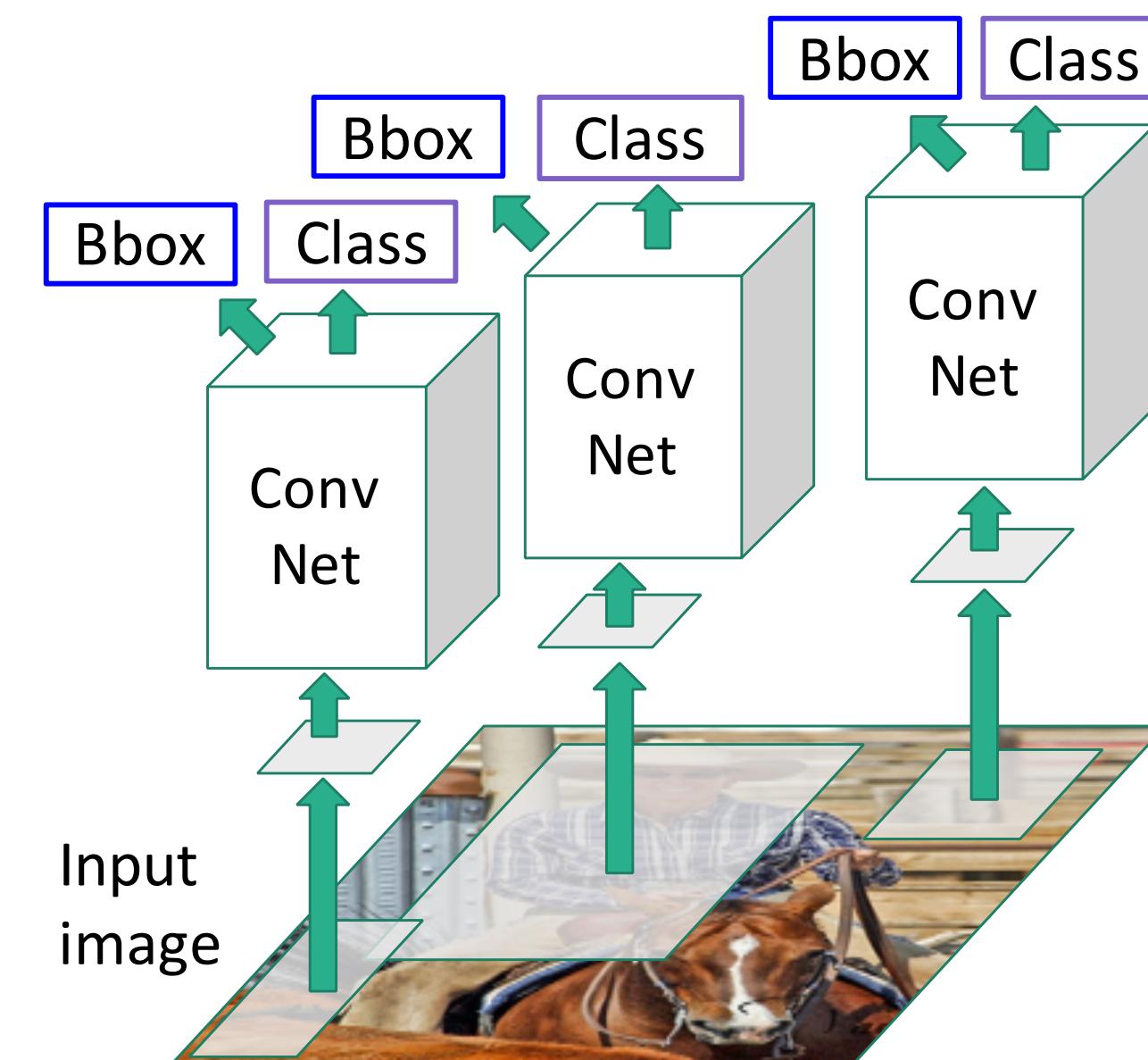
Fast R-CNN

Regions of Interest (Rois) from a proposal method

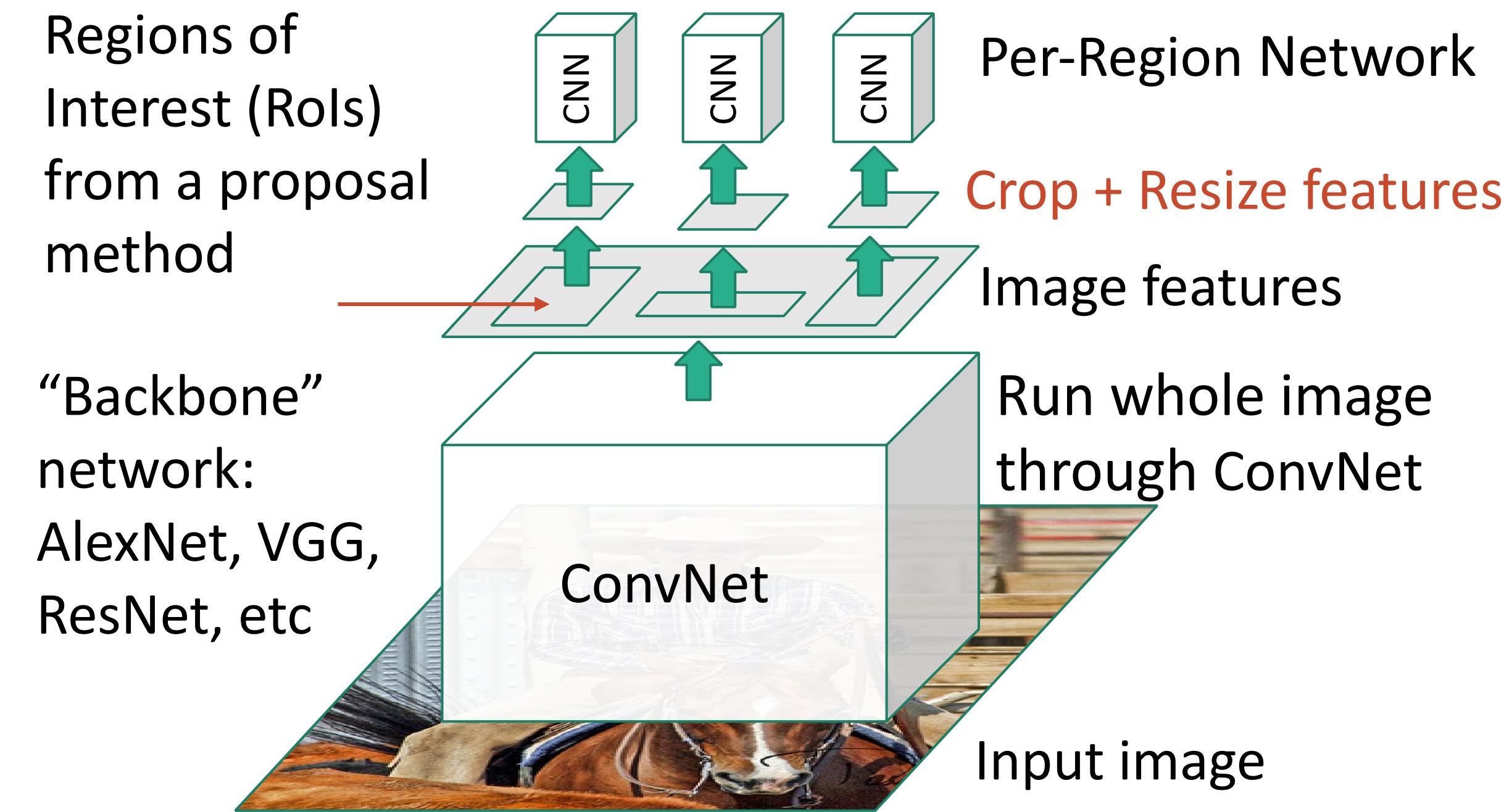
“Backbone” network:
AlexNet, VGG,
ResNet, etc



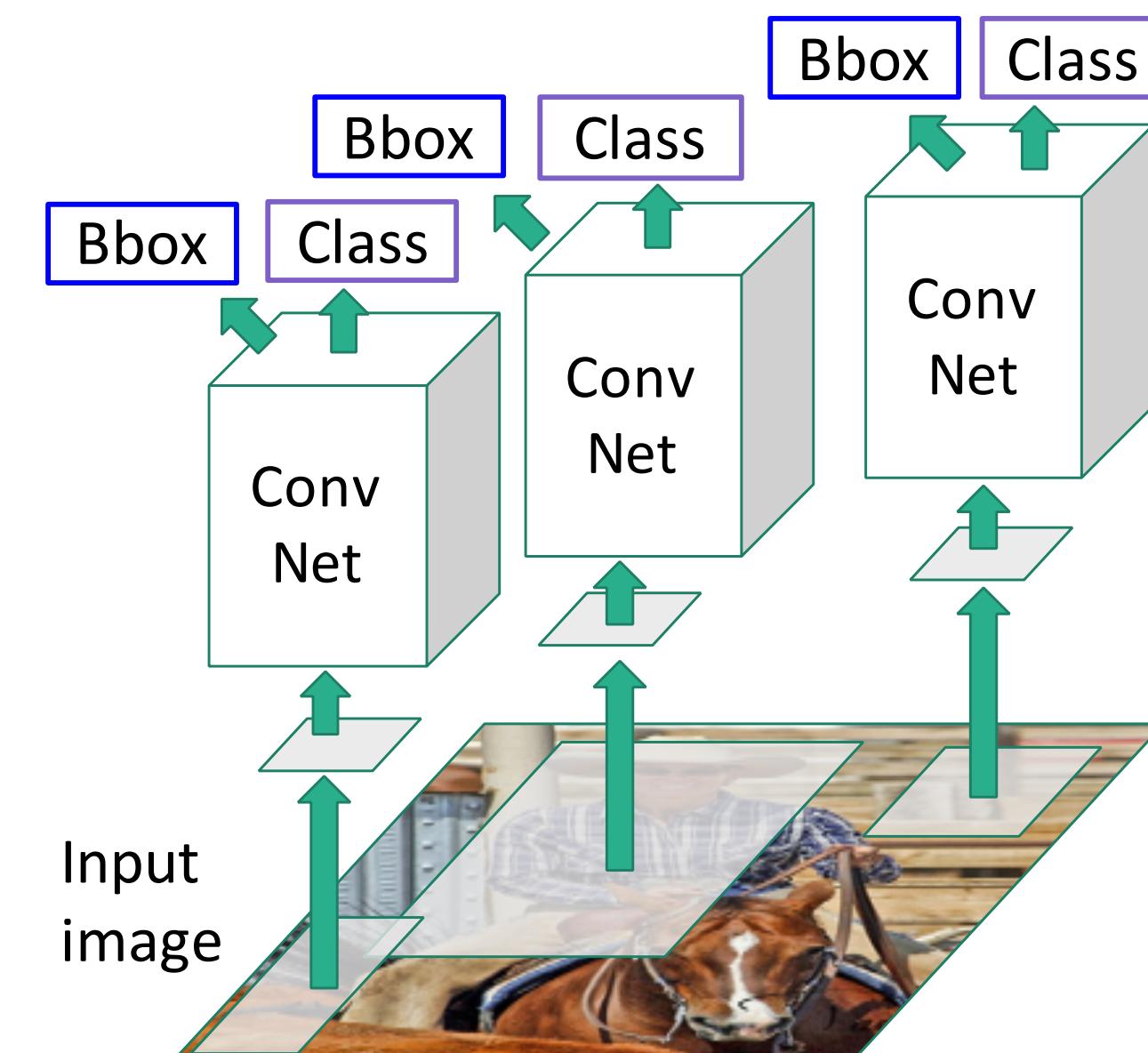
“Slow” R-CNN
Process each region independently



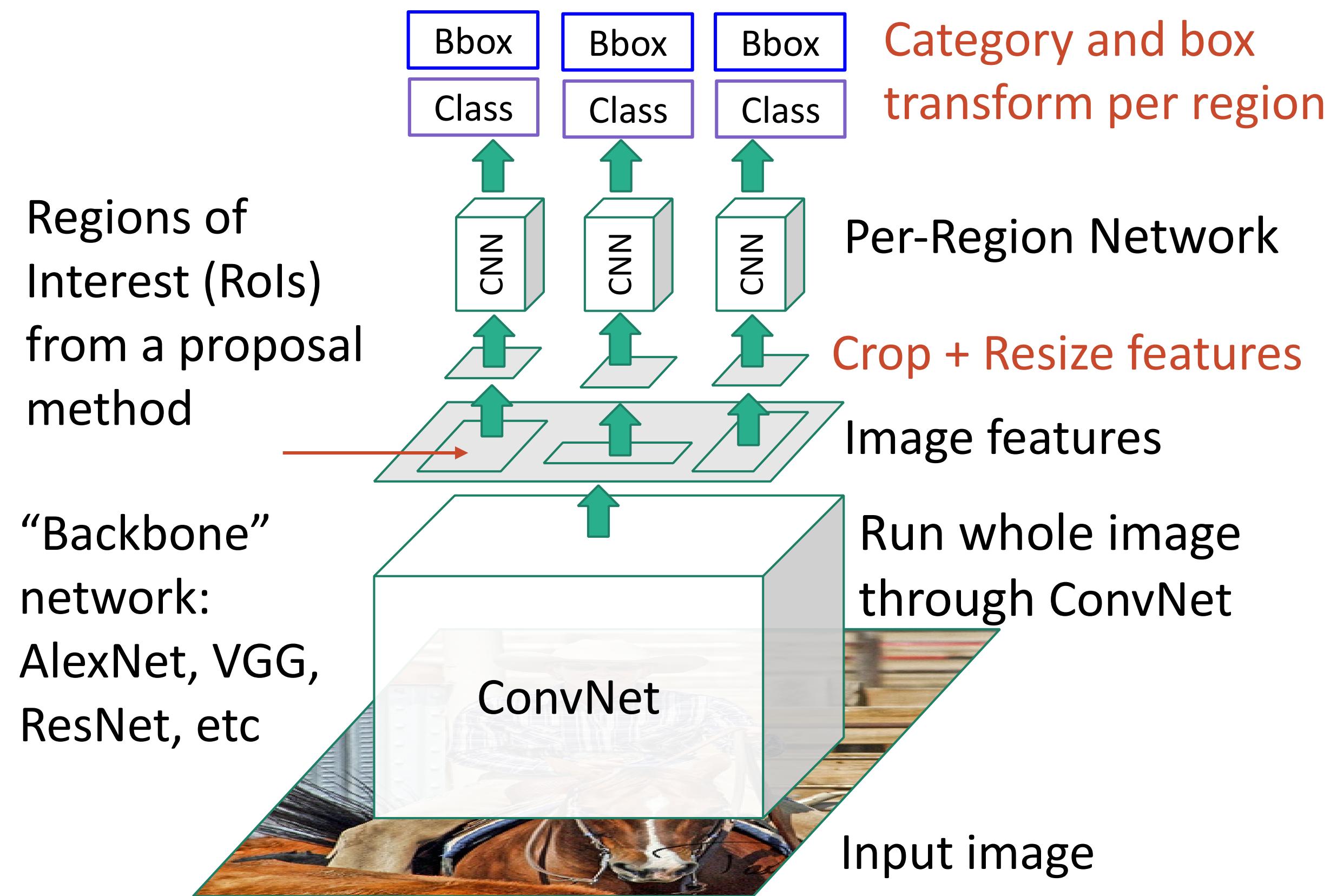
Fast R-CNN



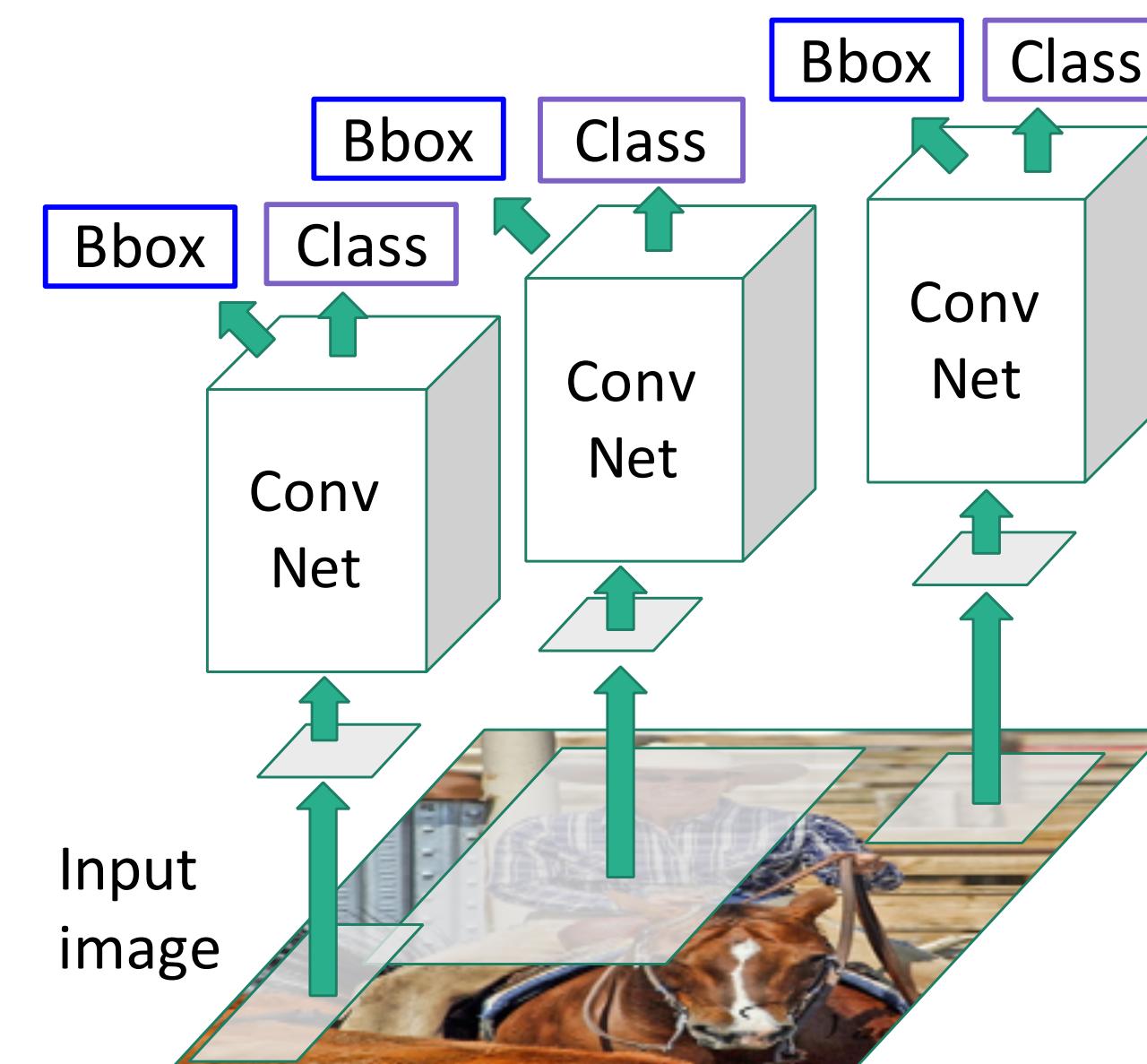
“Slow” R-CNN
Process each region
independently



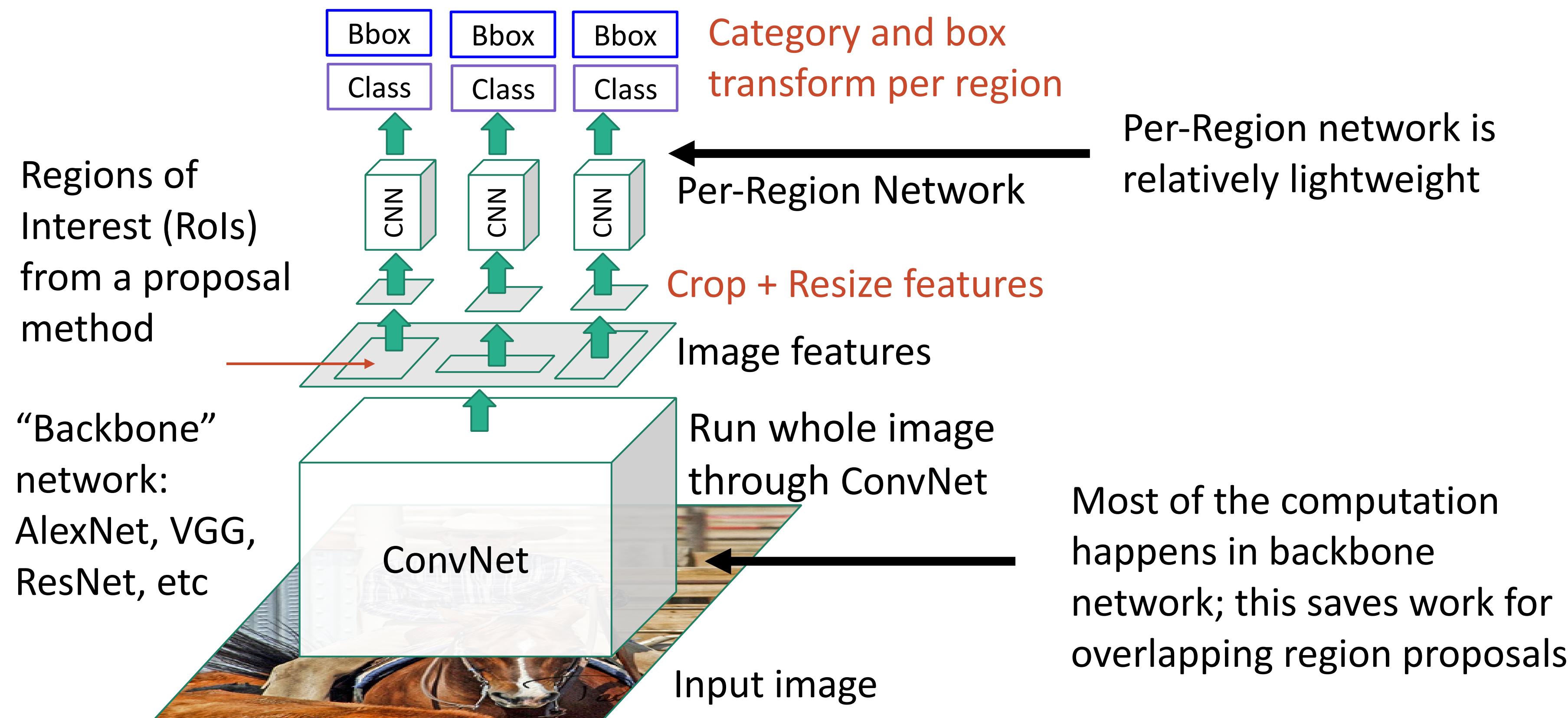
Fast R-CNN



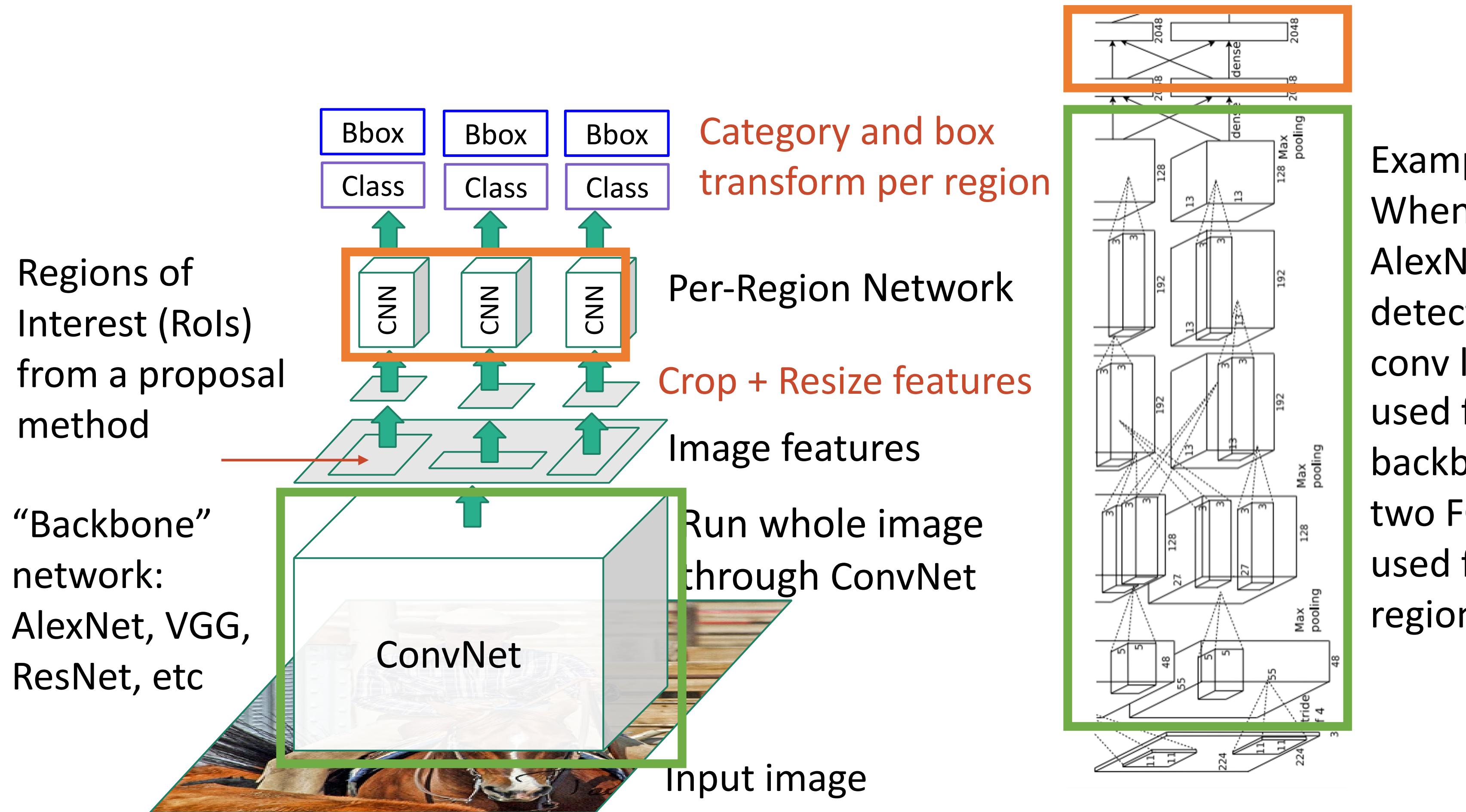
“Slow” R-CNN
Process each region independently



Fast R-CNN

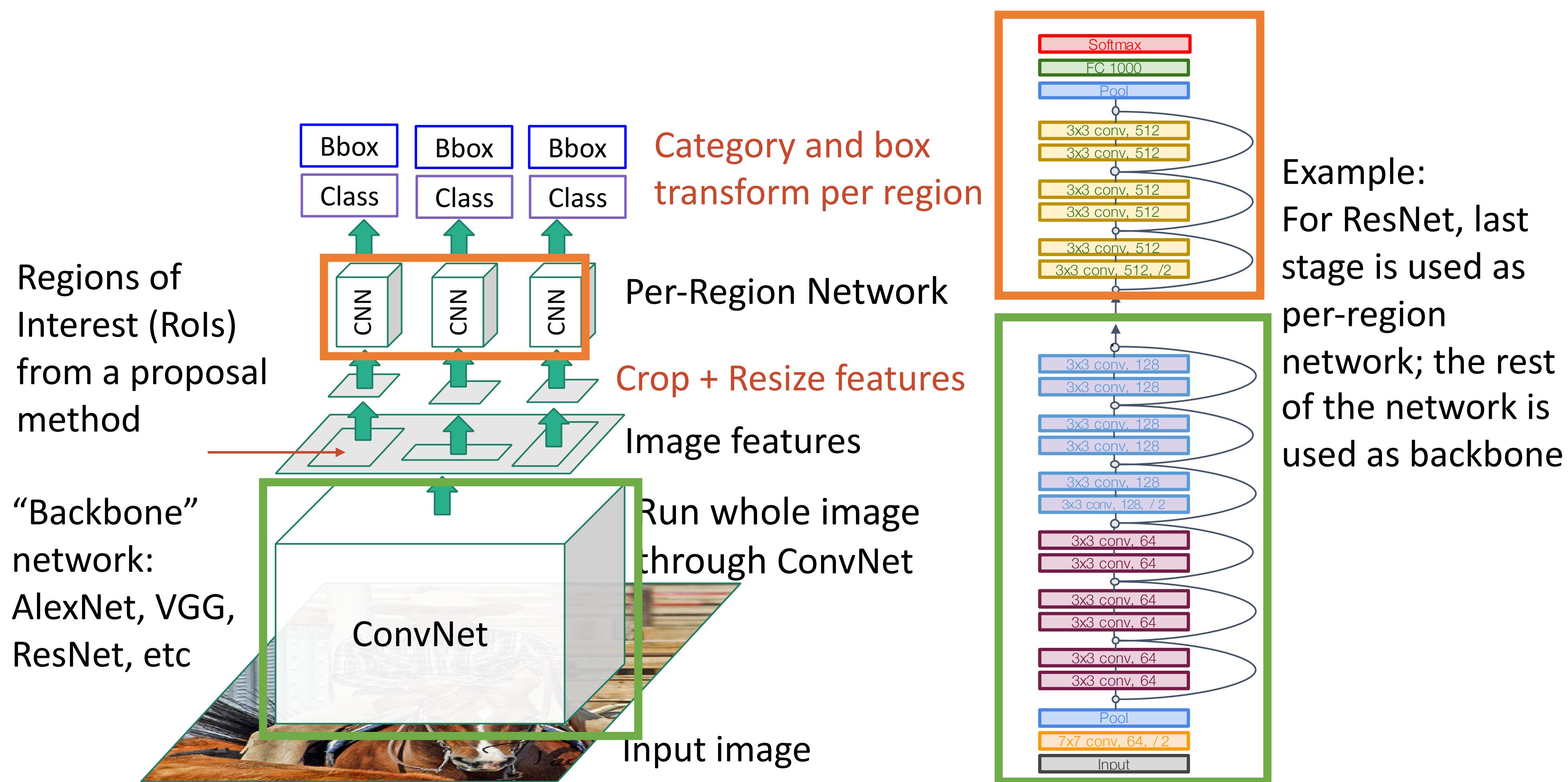


Fast R-CNN

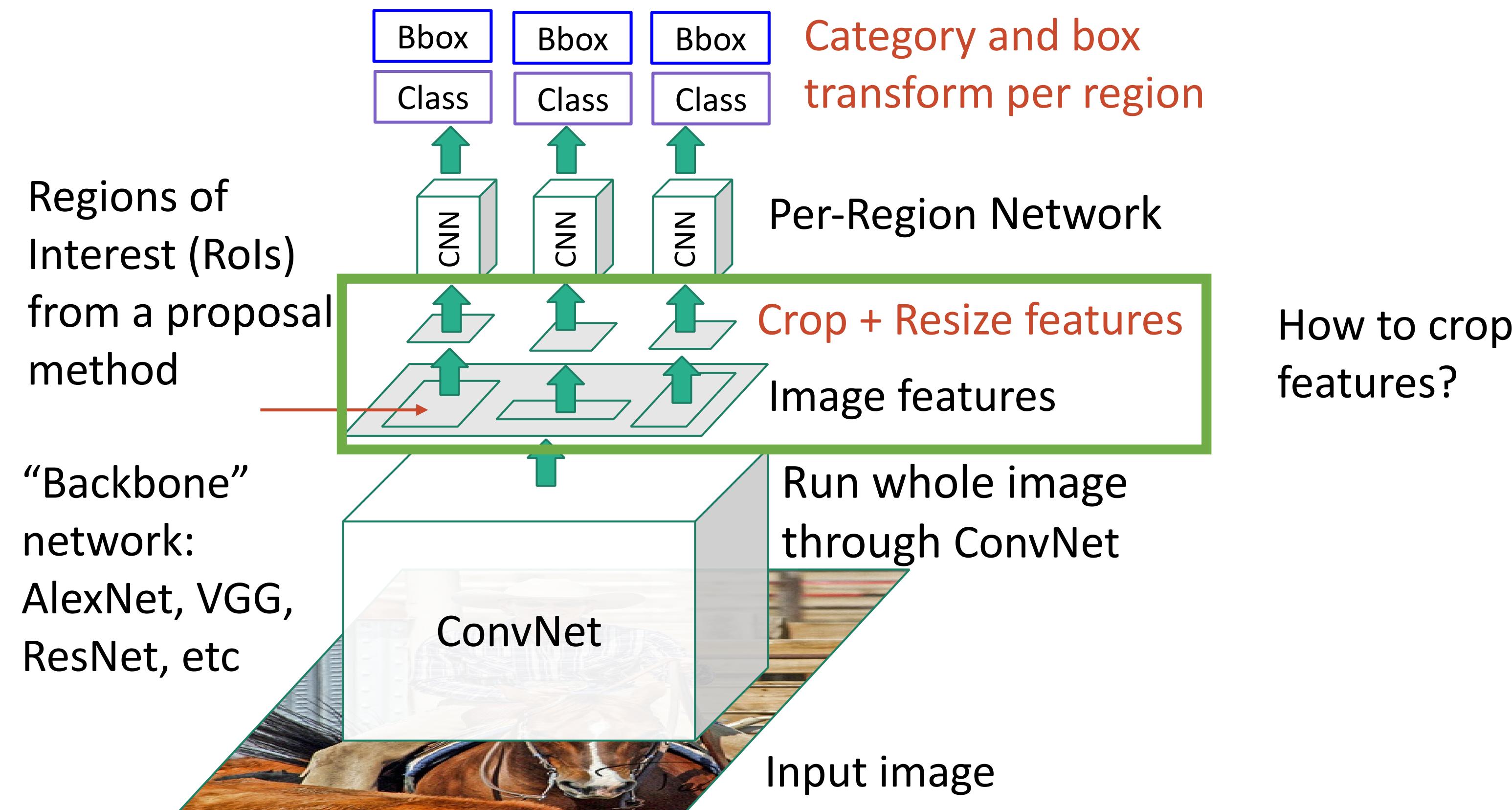


Example:
When using
AlexNet for
detection, five
conv layers are
used for
backbone and
two FC layers are
used for per-
region network

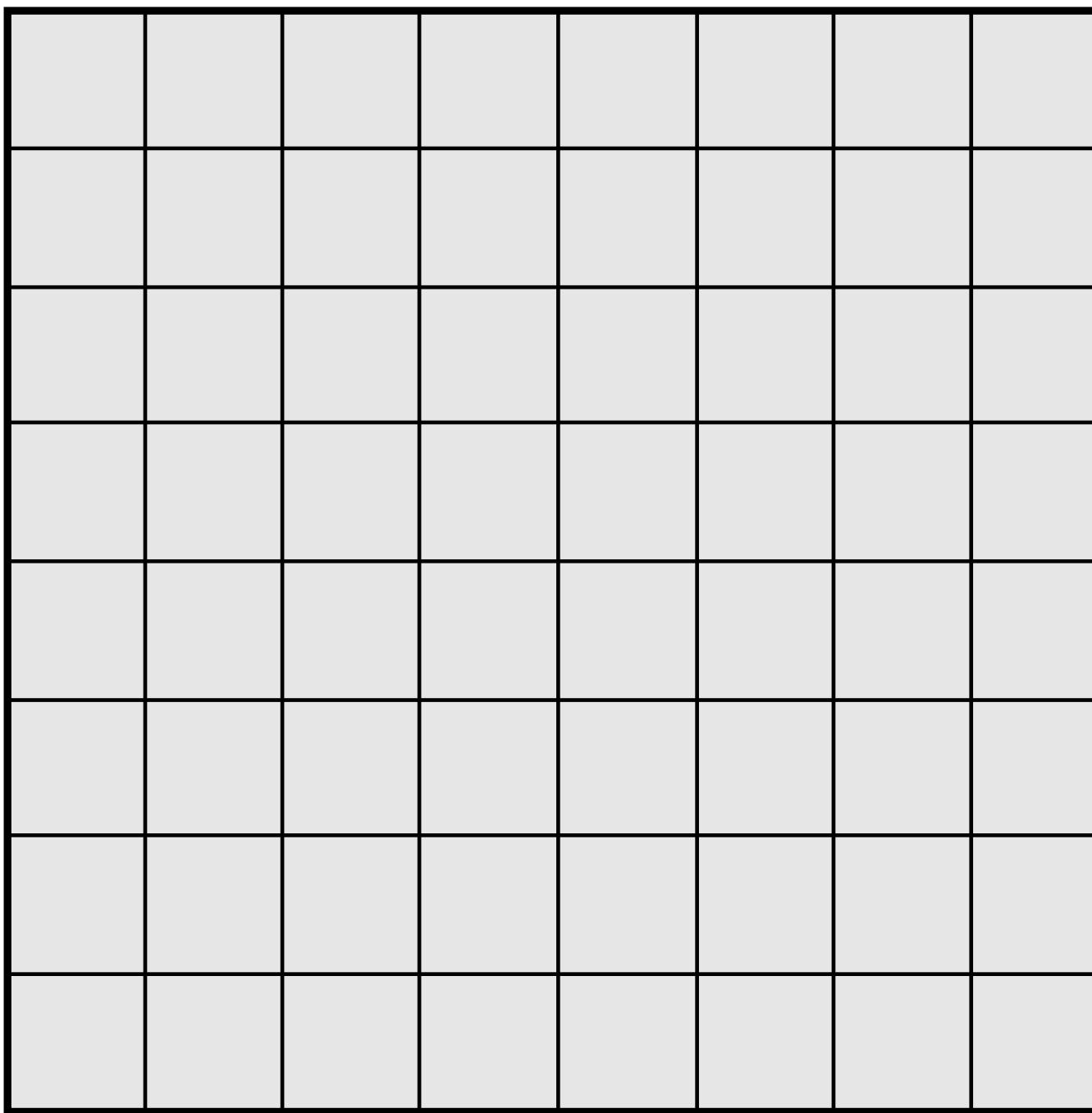
Fast R-CNN



Fast R-CNN



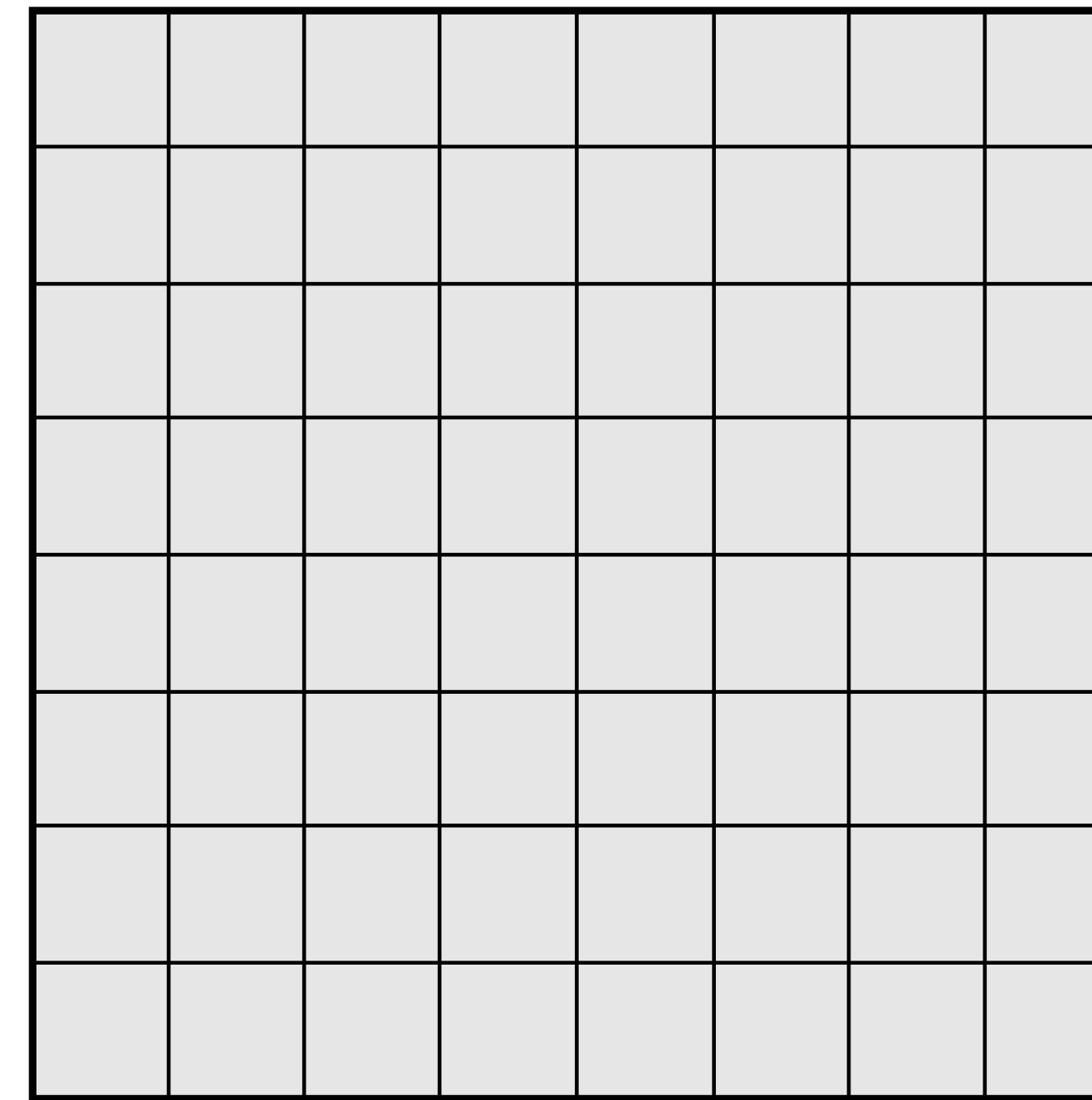
Recall: Receptive Fields



Input Image: 8×8

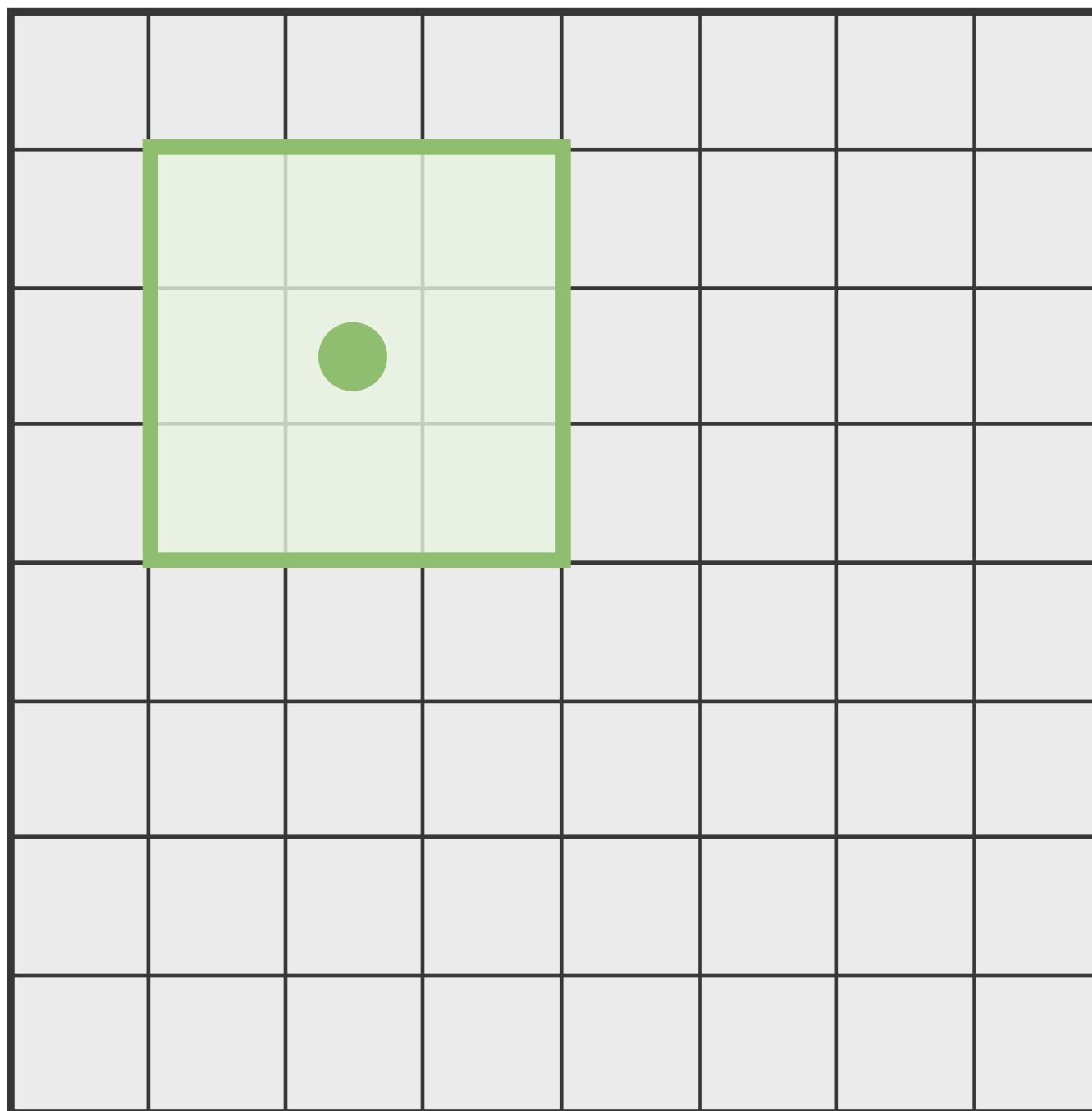
Every position in the output feature map depends on a 3×3 receptive field in the input

3x3 Conv
Stride 1, pad 1



Output Image: 8×8

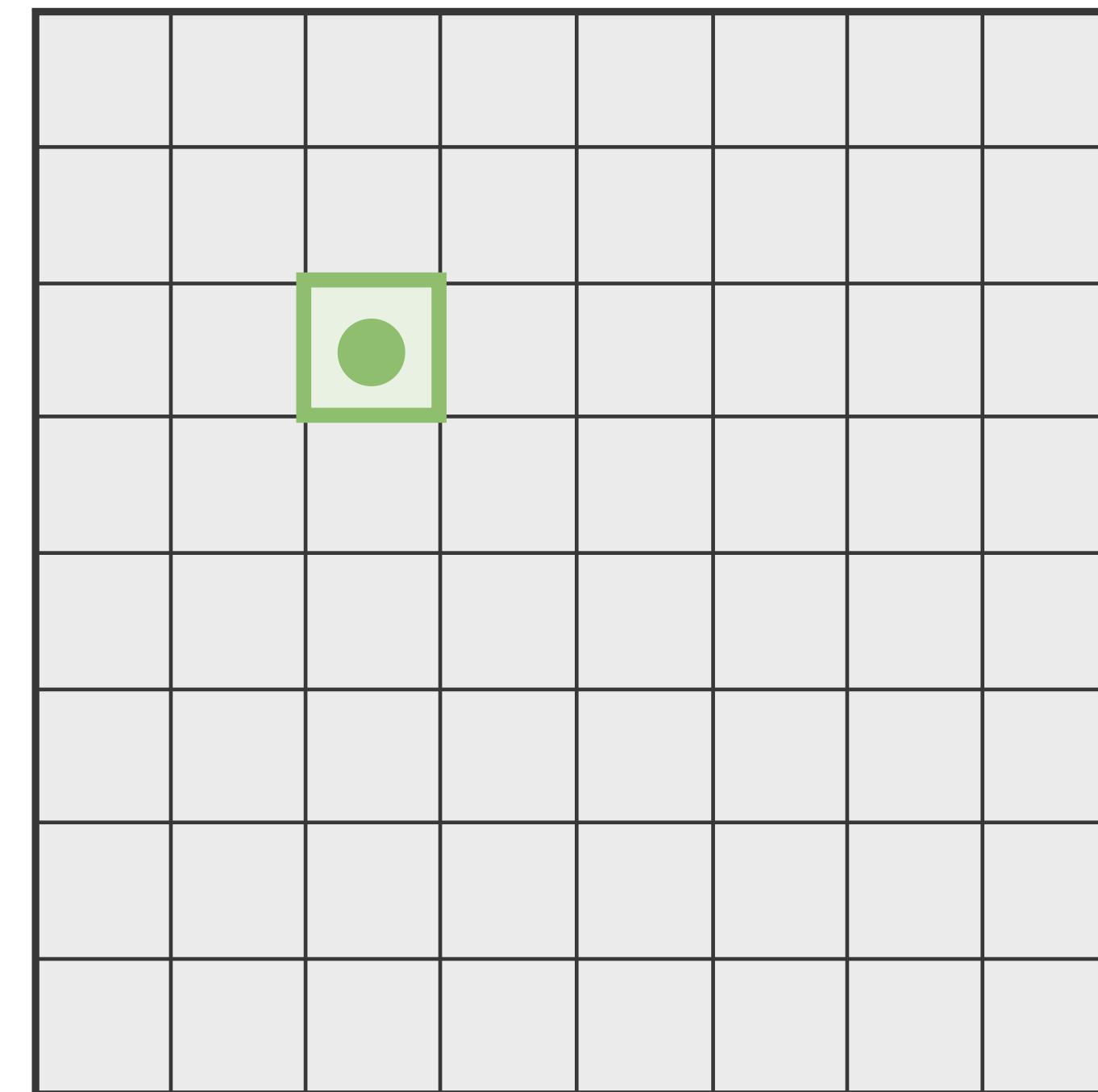
Recall: Receptive Fields



Input Image: 8 x 8

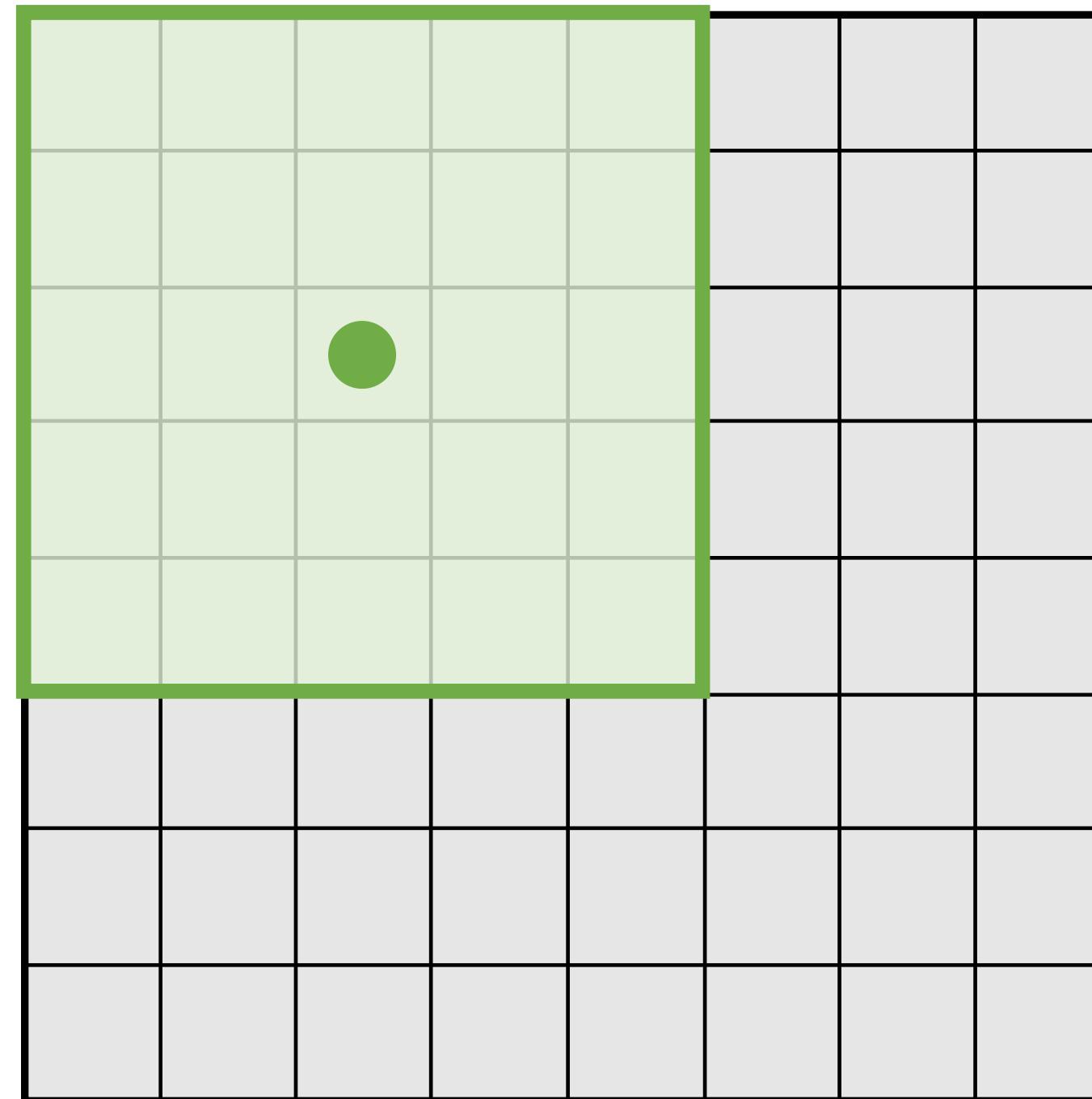
Every position in the output feature map depends on a 3x3 receptive field in the input

3x3 Conv
Stride 1, pad 1



Output Image: 8 x 8

Recall: Receptive Fields

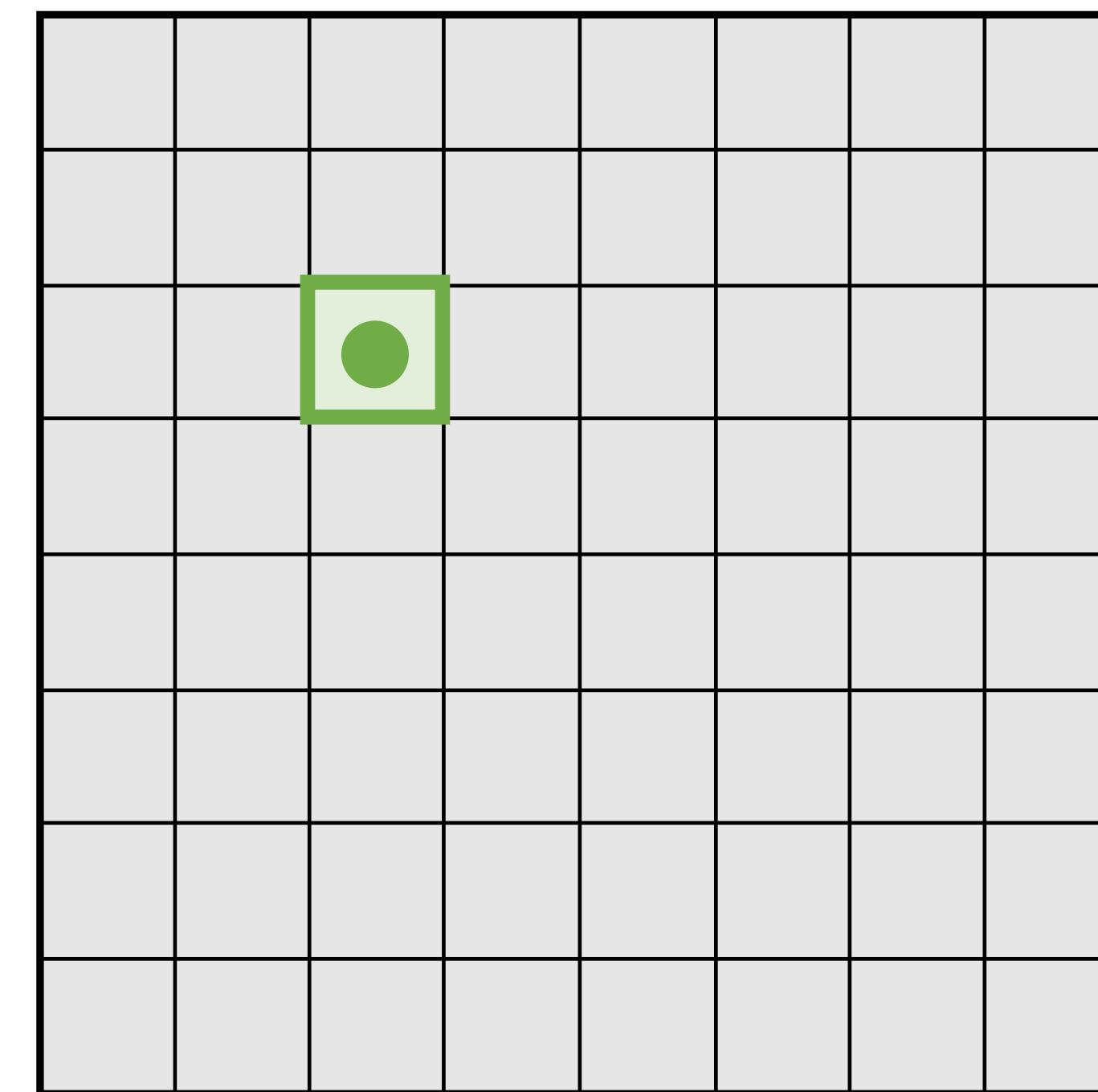


Input Image: 8 x 8

Every position in the output feature map depends on a 5x5 receptive field in the input

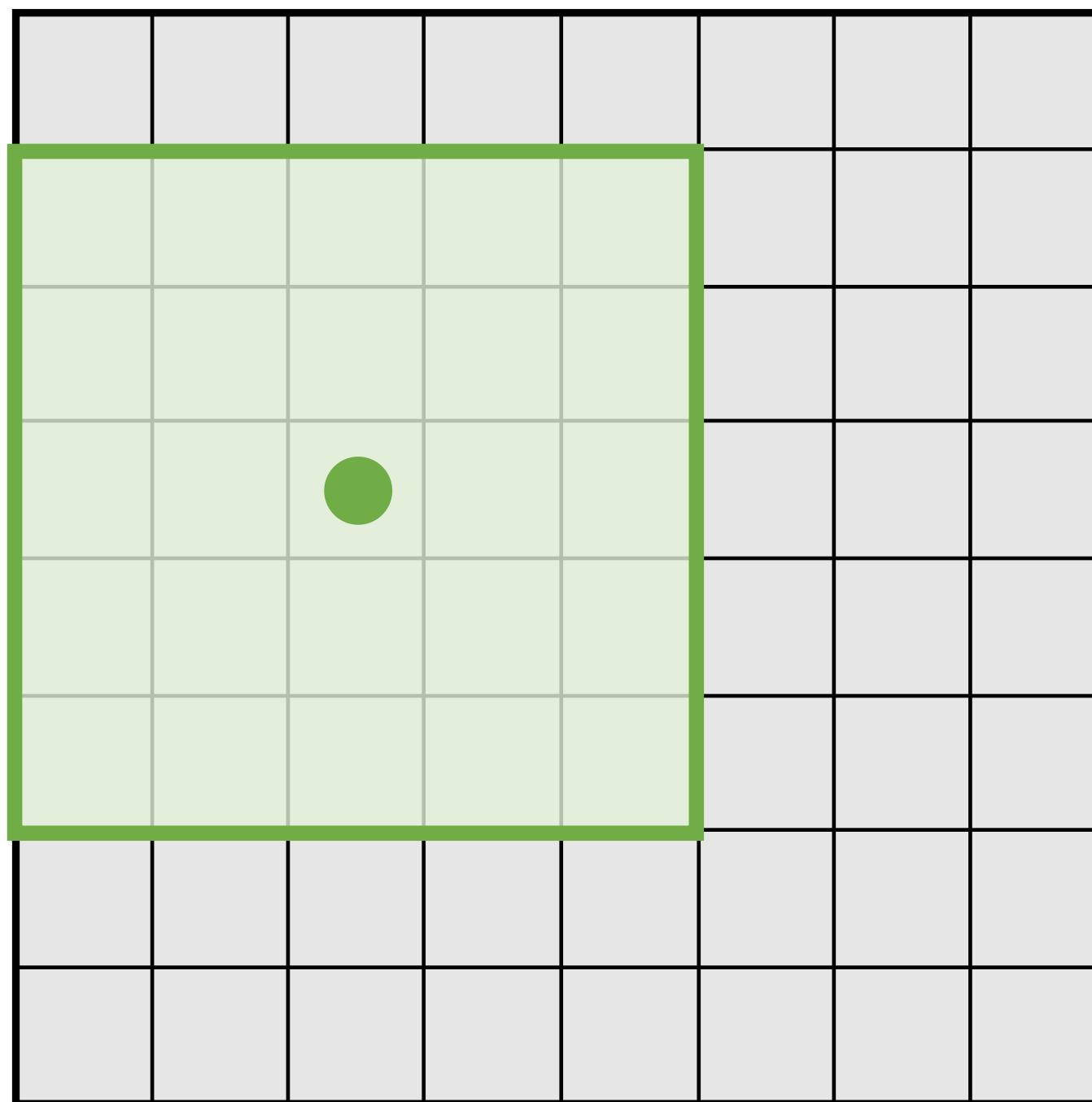
3x3 Conv
Stride 1, pad 1

3x3 Conv
Stride 1, pad 1



Output Image: 8 x 8

Recall: Receptive Fields

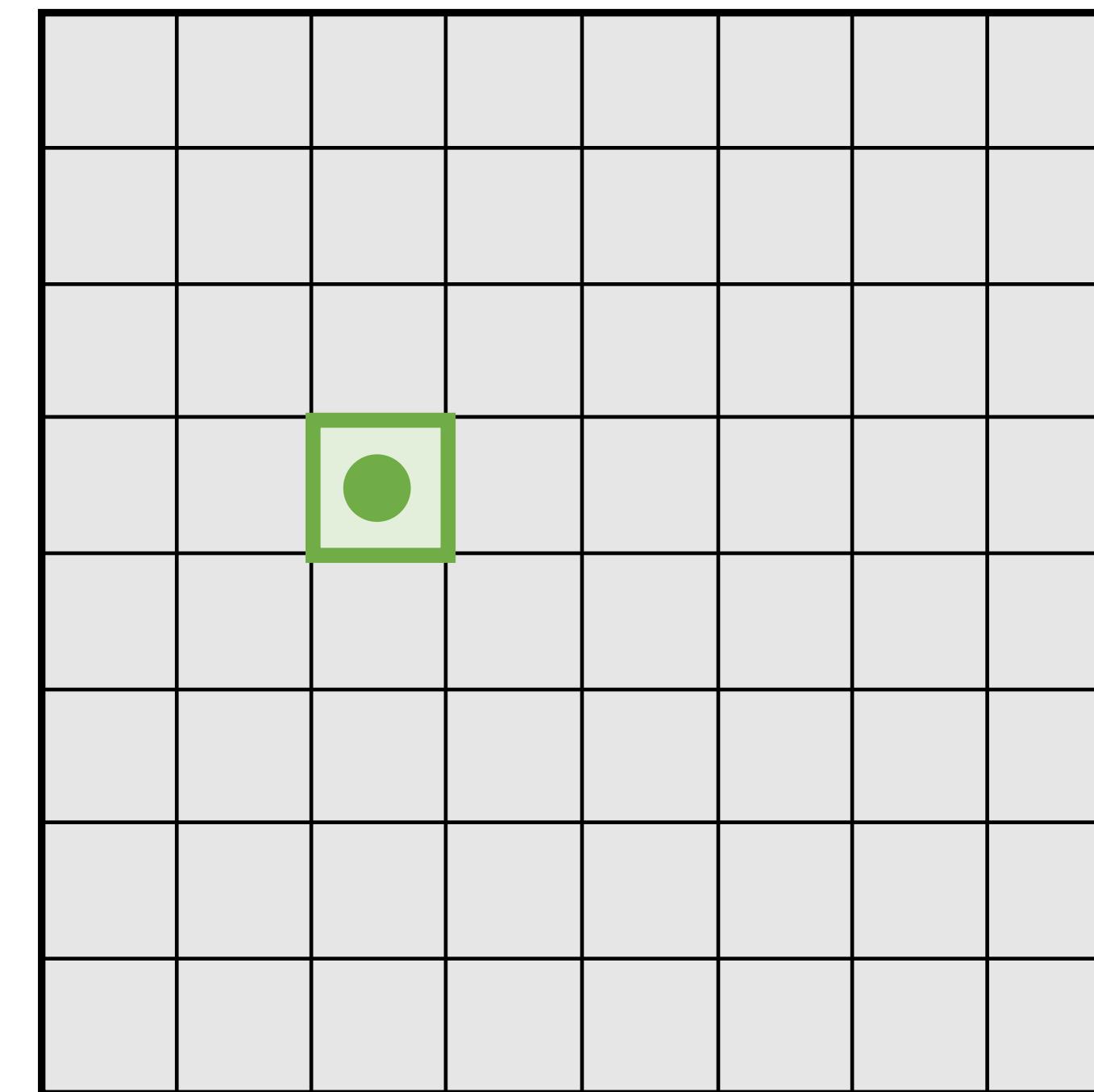


Input Image: 8 x 8

Moving one unit in the output space also moves the receptive field by one

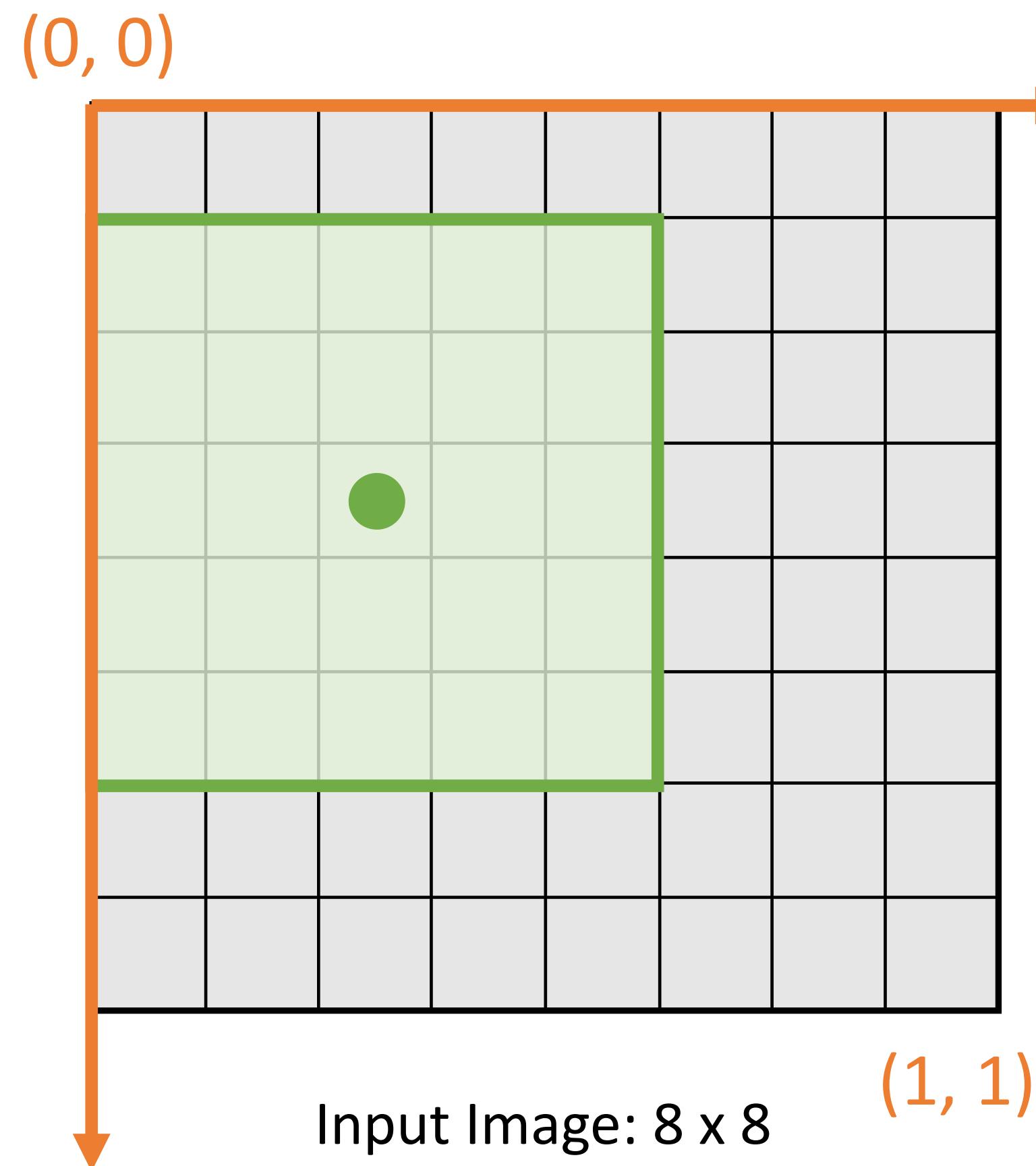
3x3 Conv
Stride 1, pad 1

3x3 Conv
Stride 1, pad 1



Output Image: 8 x 8

Recall: Receptive Fields

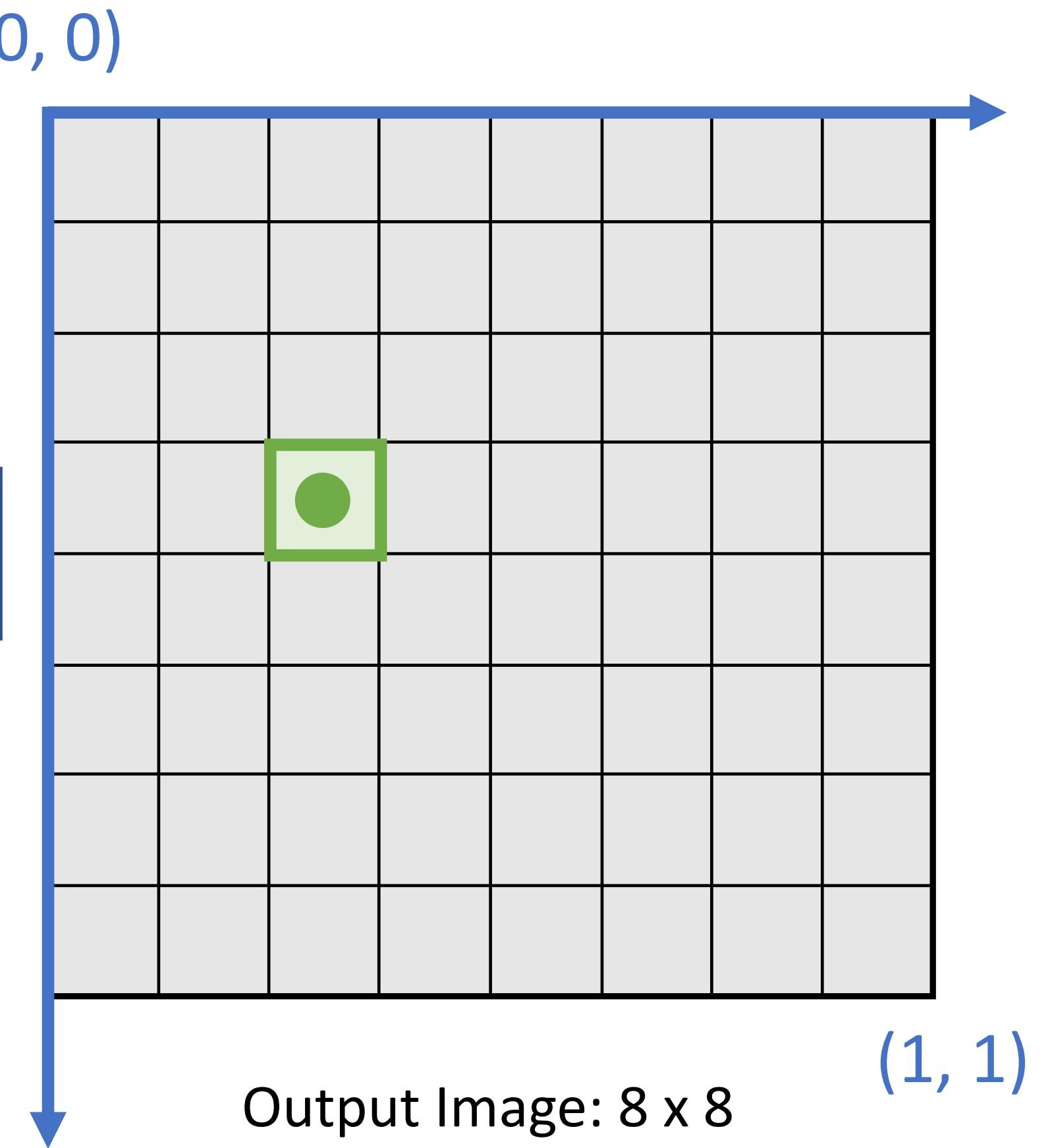


Moving one unit in the output space also moves the receptive field by one

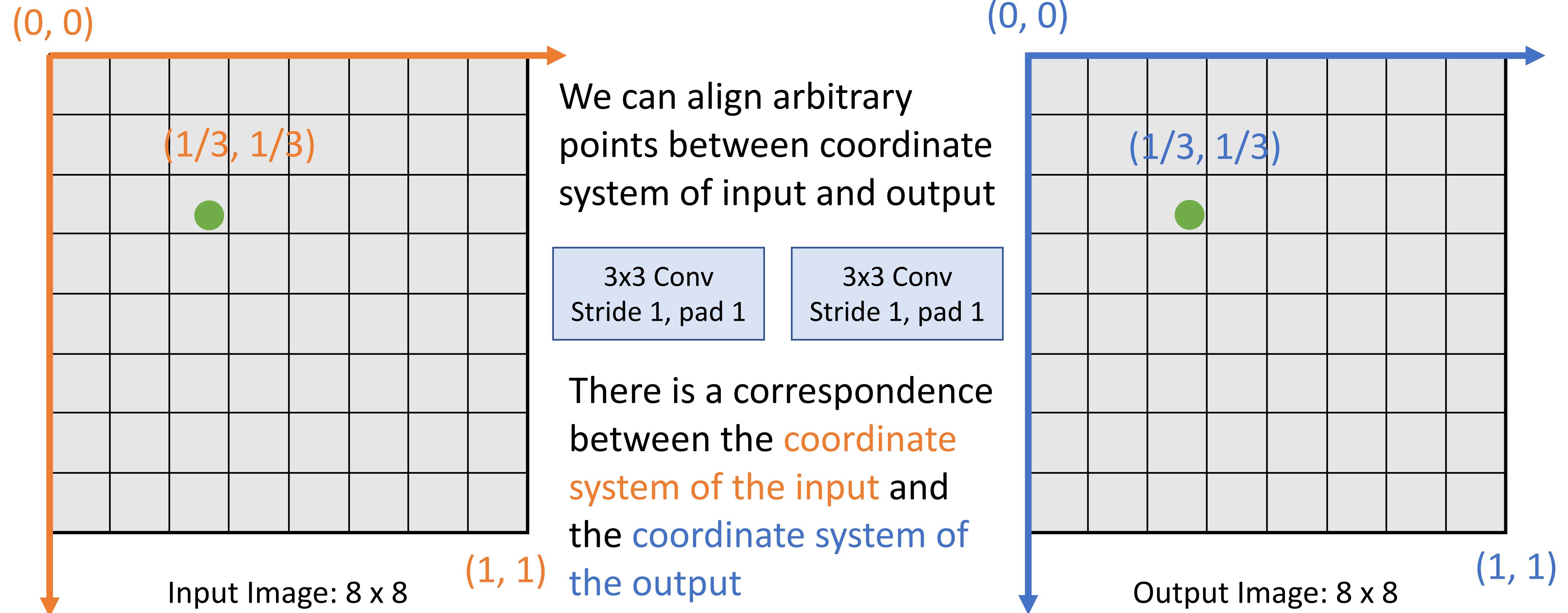
3x3 Conv
Stride 1, pad 1

3x3 Conv
Stride 1, pad 1

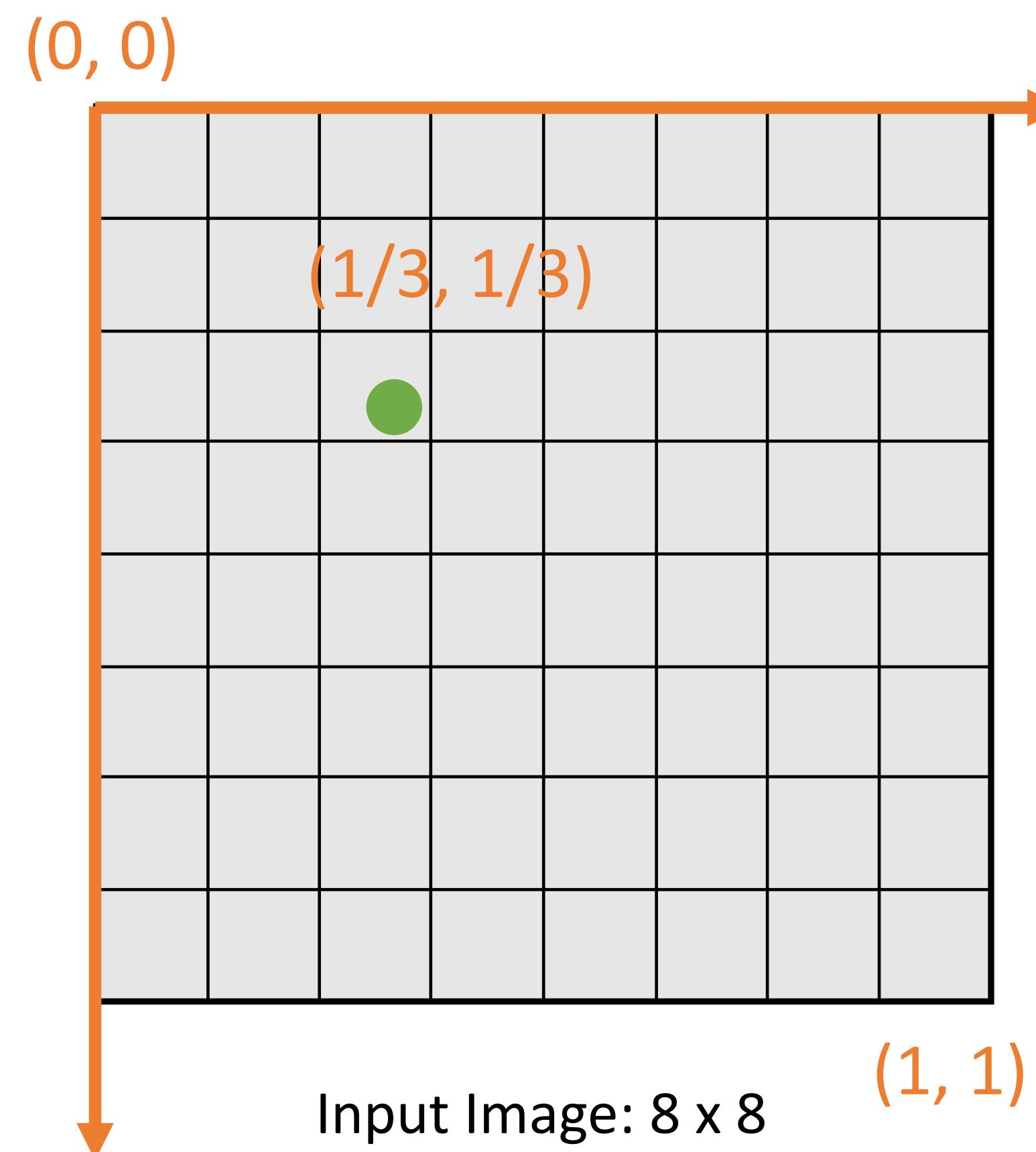
There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**



Projecting Points



Projecting Points

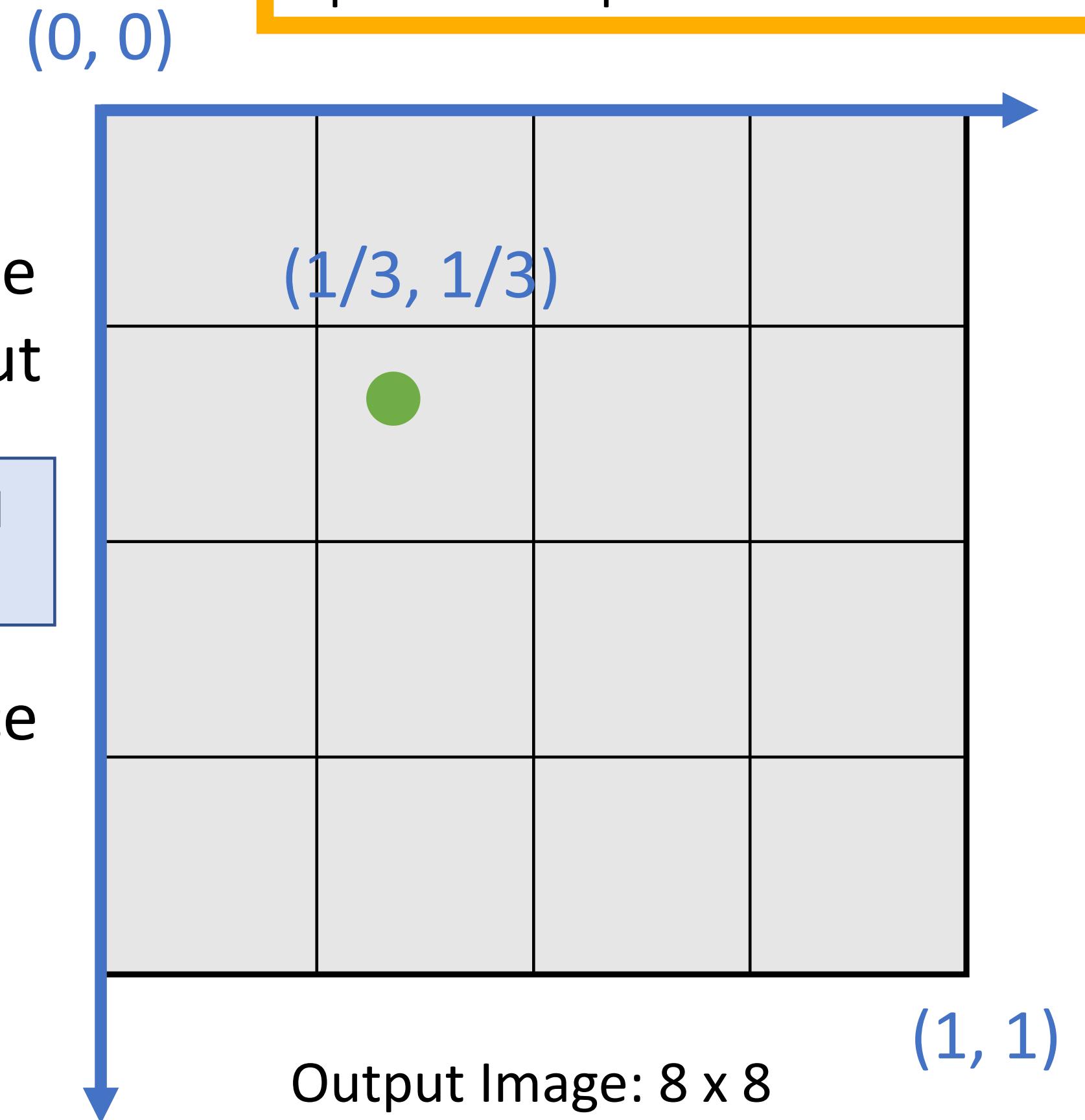


We can align arbitrary points between coordinate system of input and output

3x3 Conv
Stride 1, pad 1

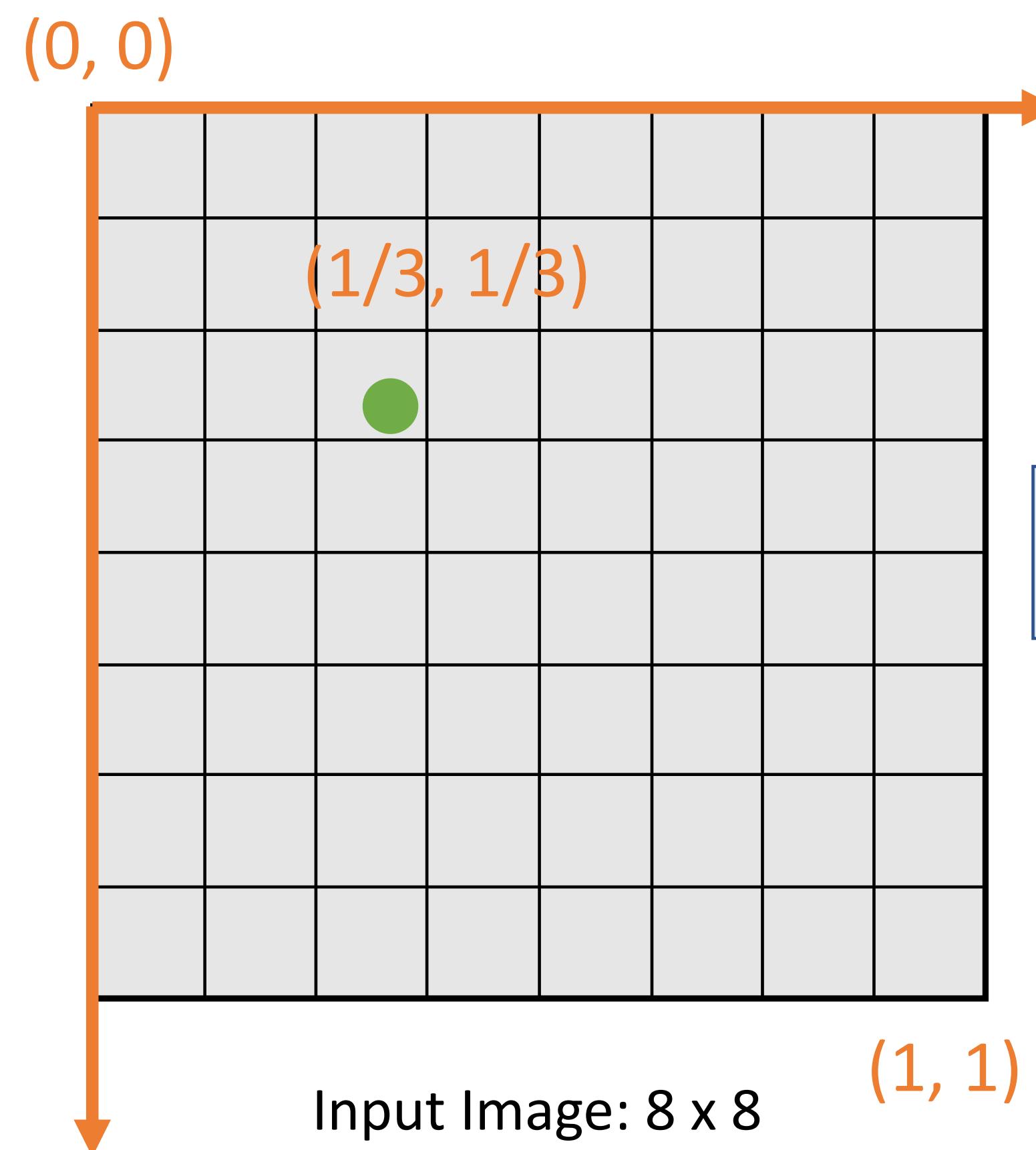
2x2 MaxPool
Stride 2

There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**



Same logic holds for more complicated CNNs, even if spatial resolution of input and output are different

Projecting Points

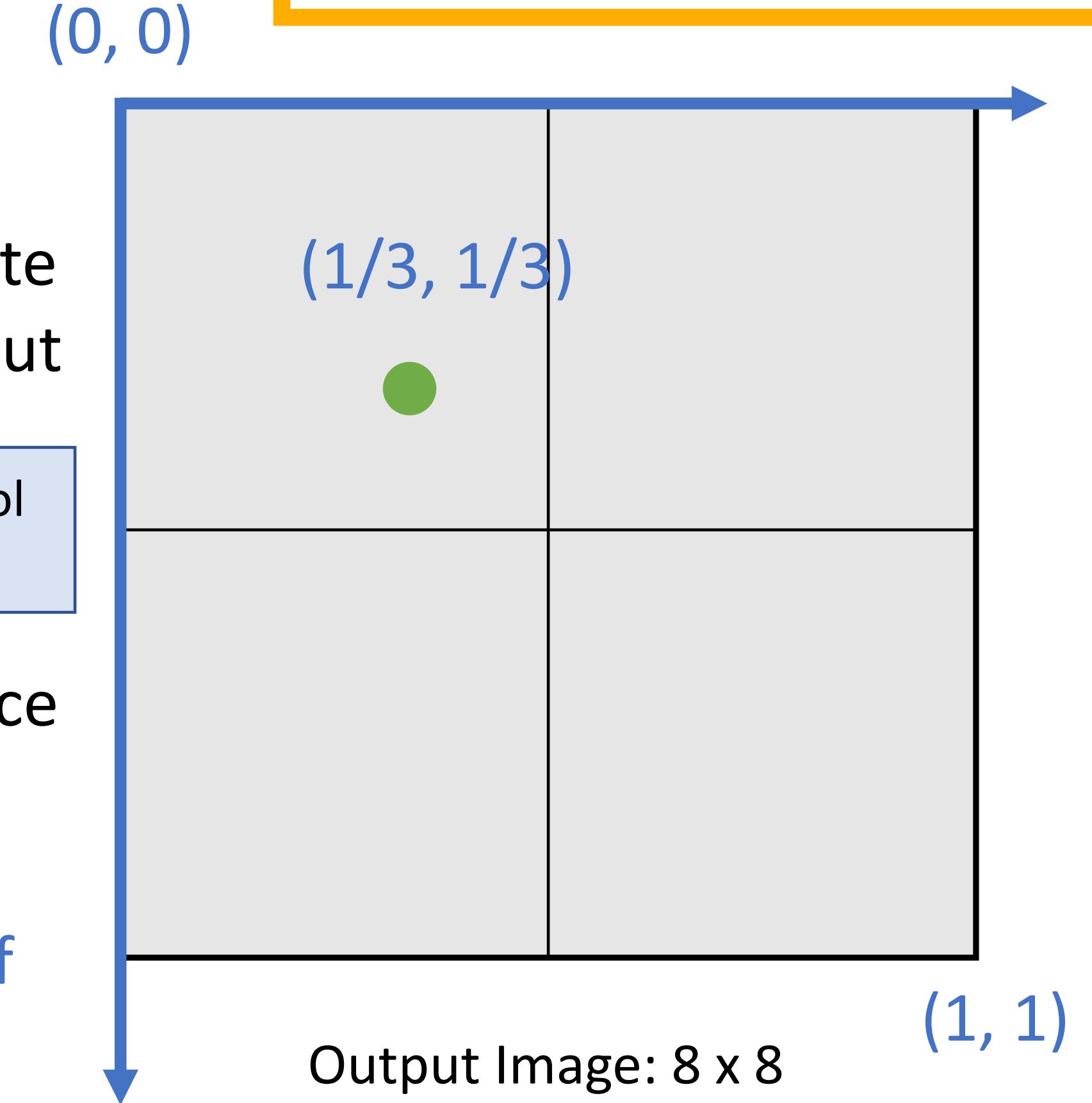


We can align arbitrary points between coordinate system of input and output

3x3 Conv
Stride 1, pad 1

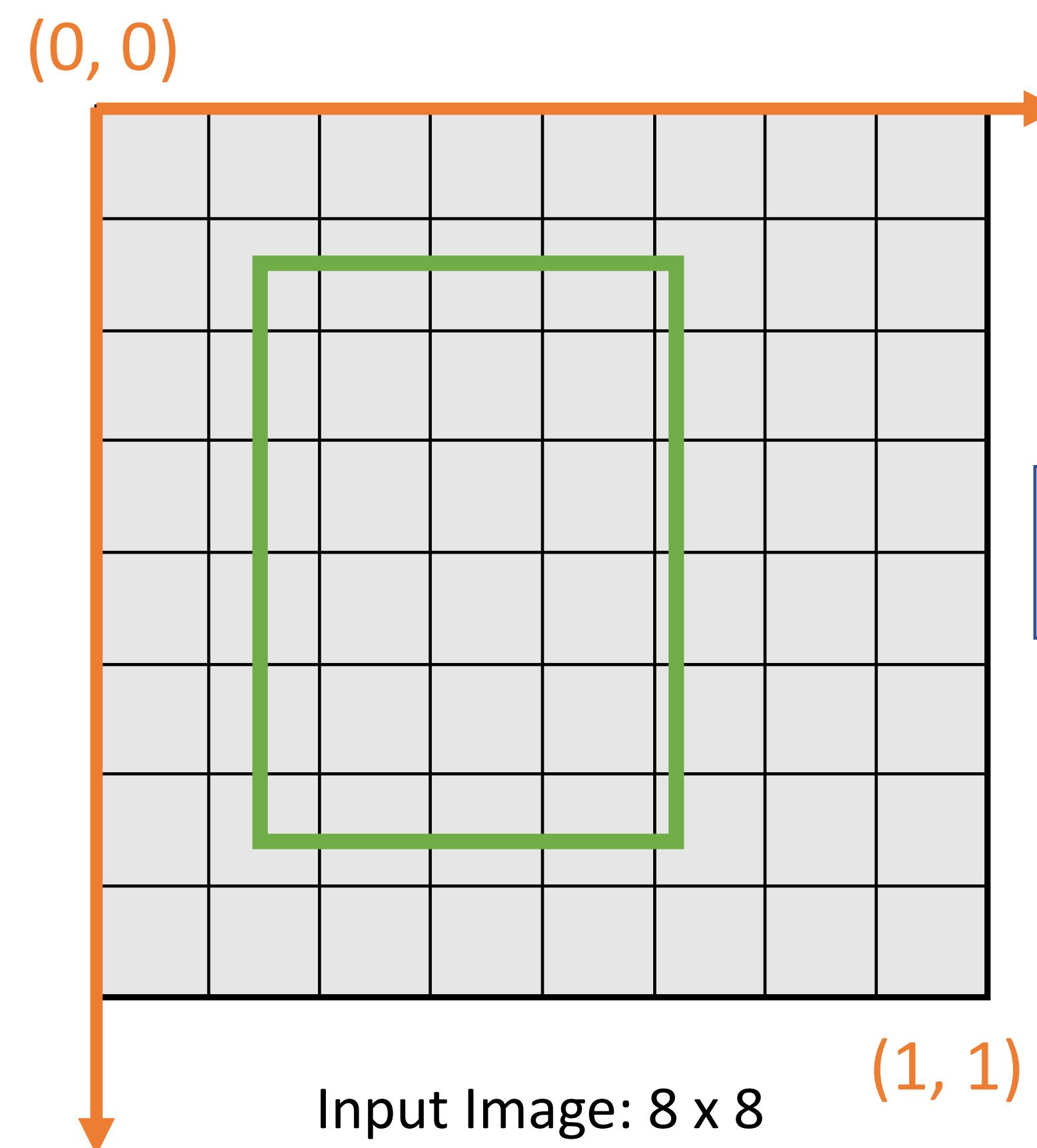
4x4 MaxPool
Stride 4

There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**



Same logic holds for more complicated CNNs, even if spatial resolution of input and output are different

Projecting Points

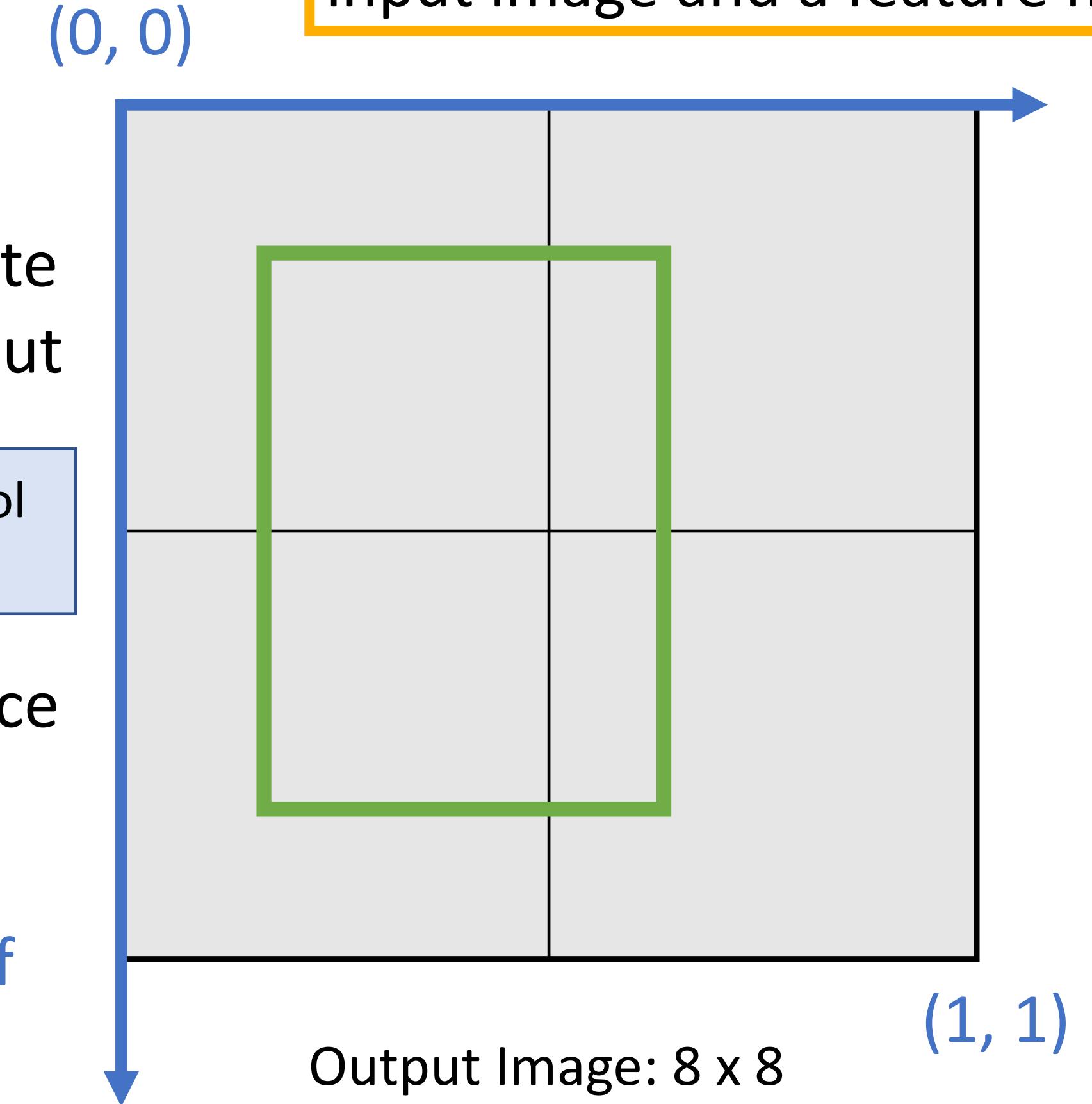


We can align arbitrary points between coordinate system of input and output

3x3 Conv
Stride 1, pad 1

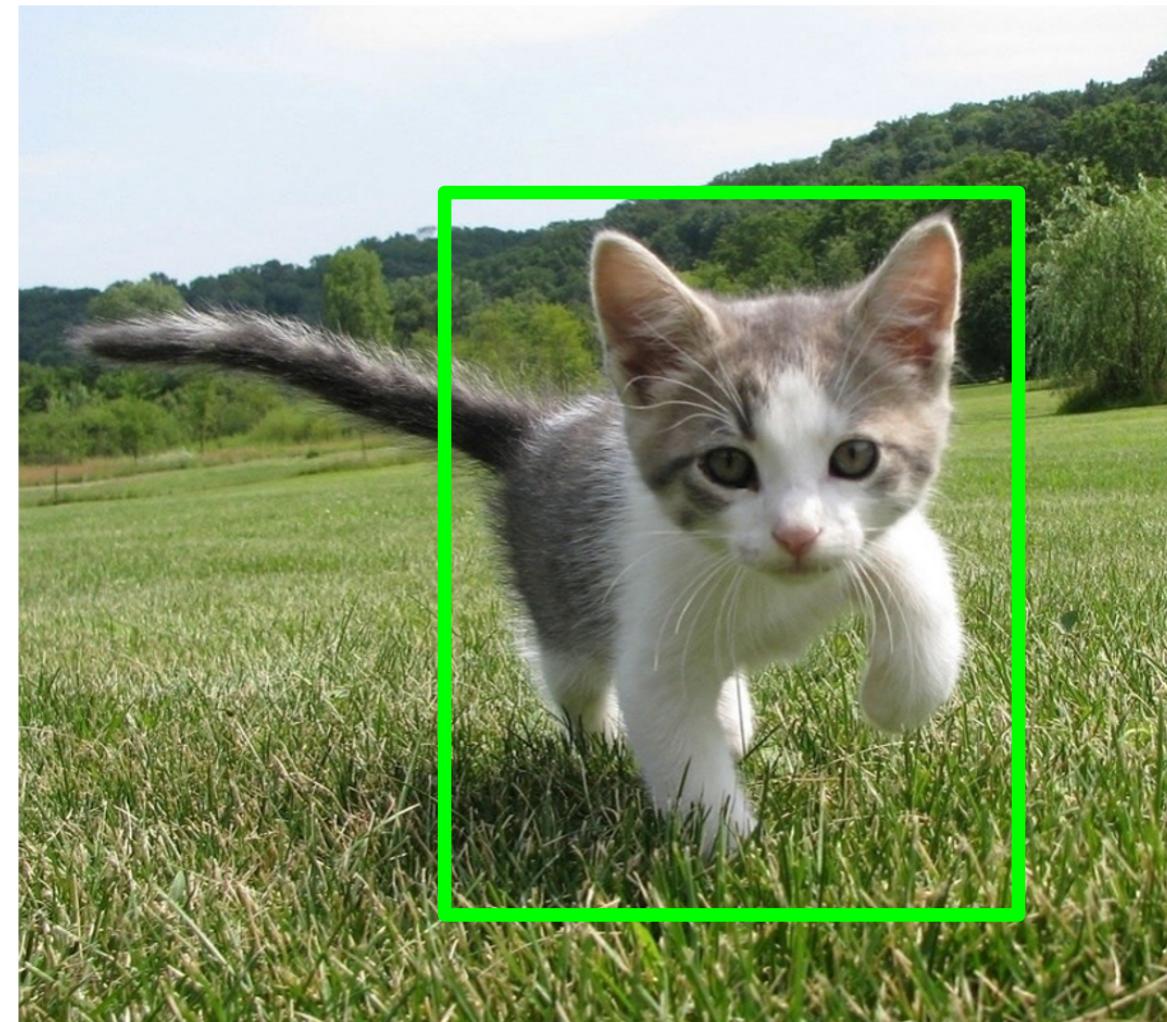
4x4 MaxPool
Stride 4

There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**



We can use this idea to project **bounding boxes** between an input image and a feature map

Cropping Features: RoI Pool



Input Image
(e.g. $3 \times 640 \times 480$)

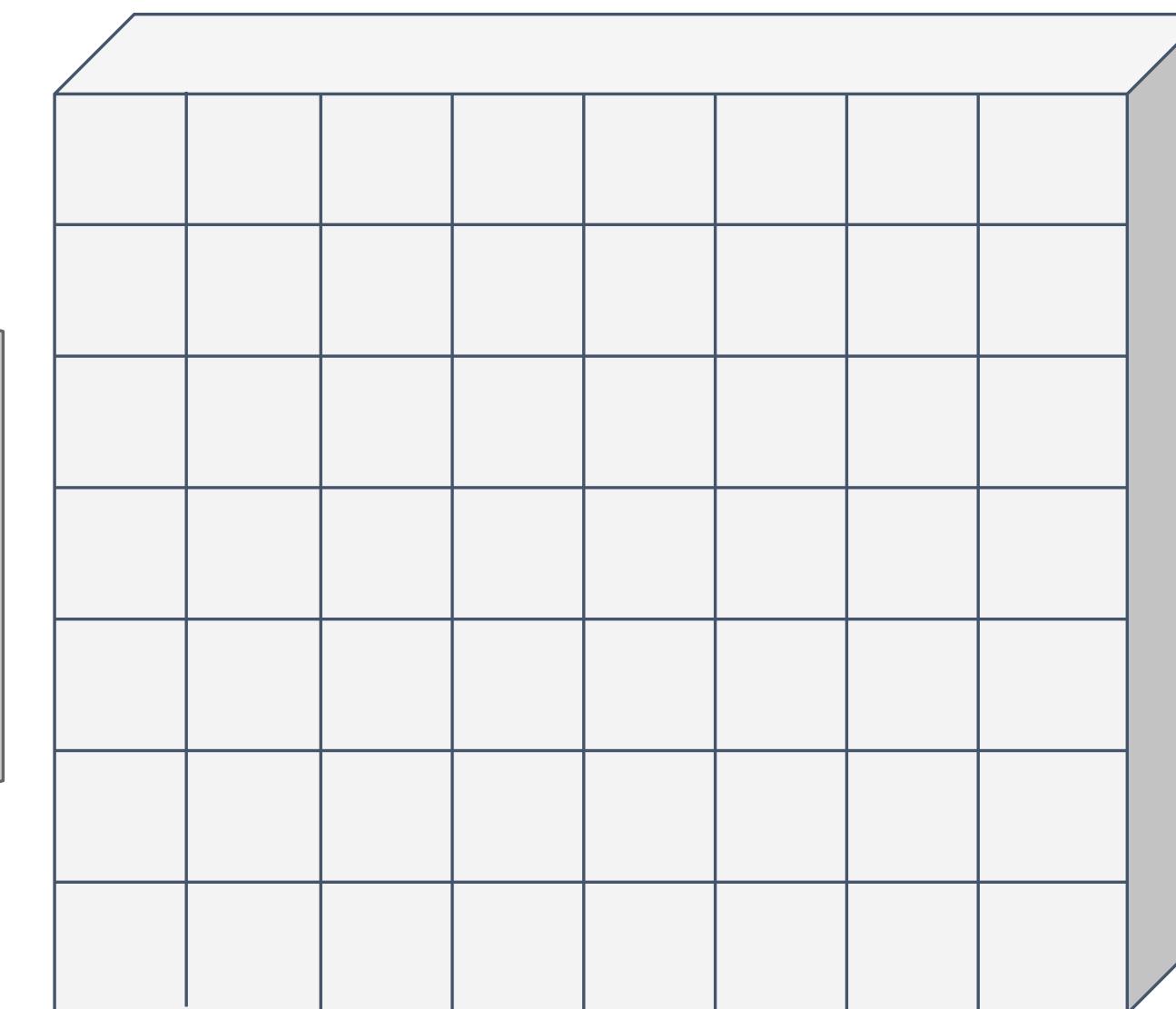
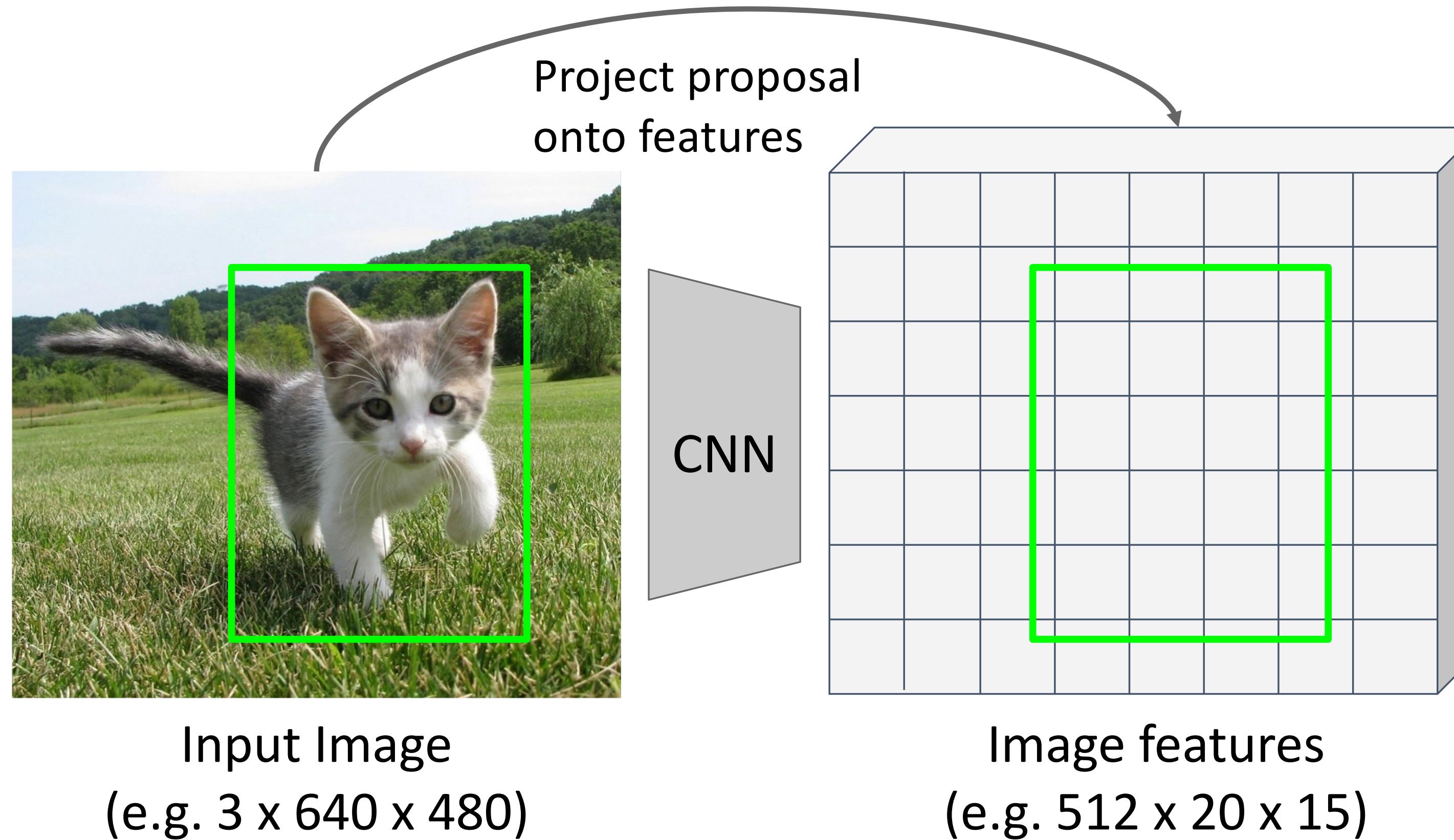


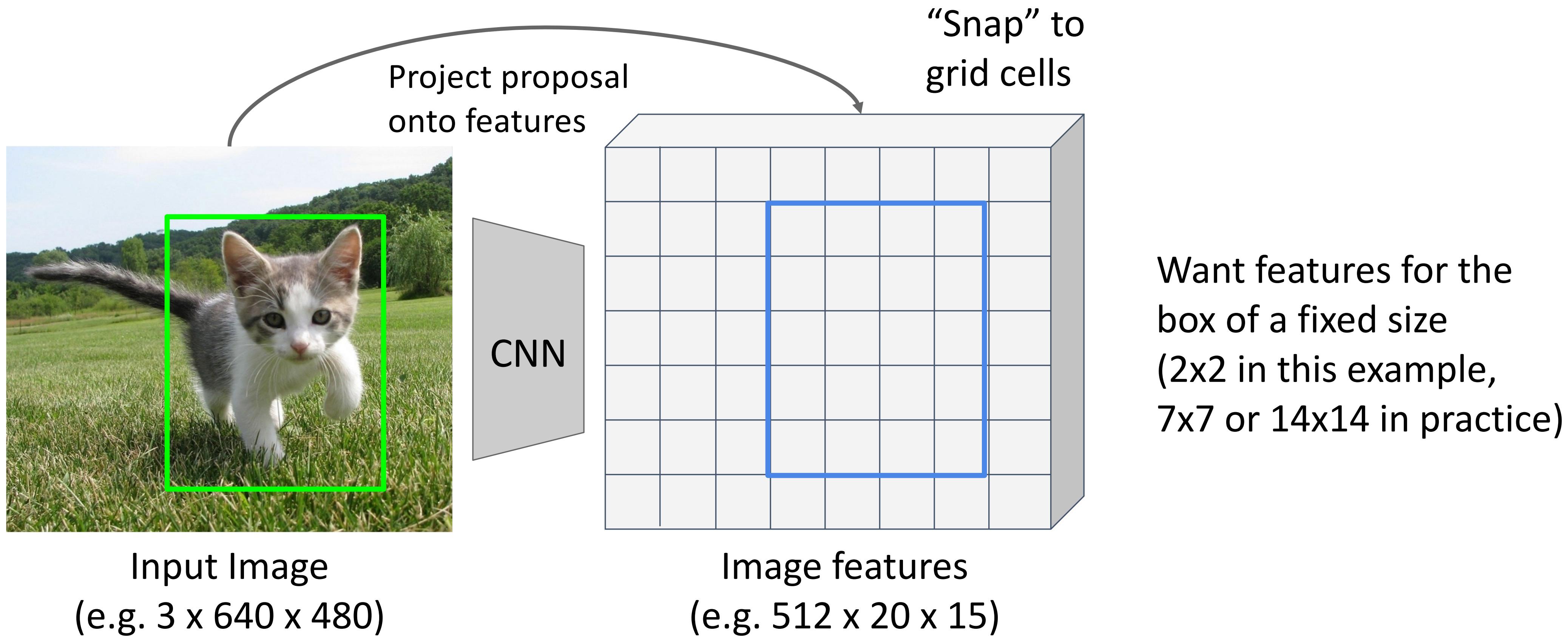
Image features
(e.g. $512 \times 20 \times 15$)

Want features for the
box of a fixed size
(2×2 in this example,
 7×7 or 14×14 in practice)

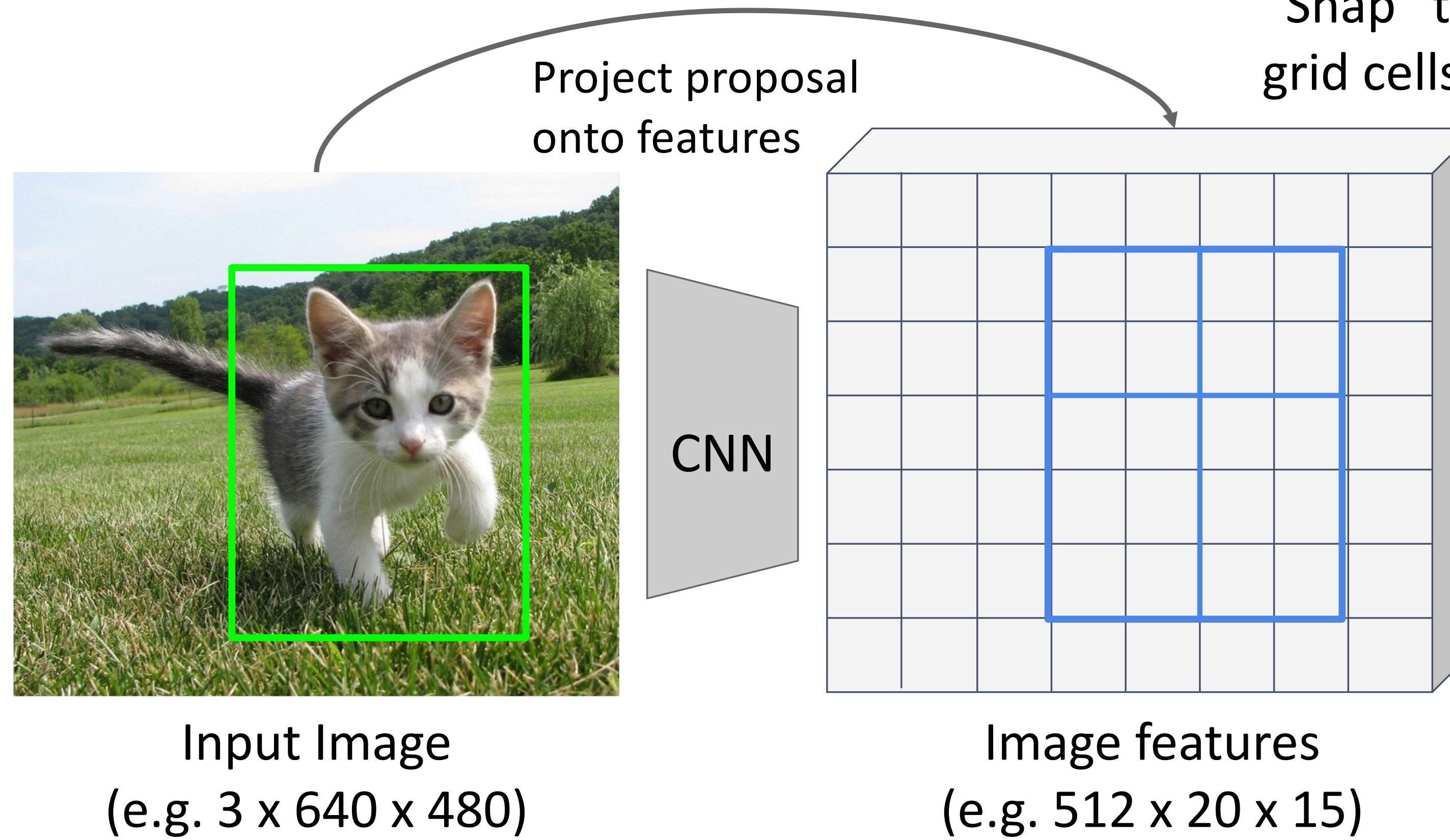
Cropping Features: RoI Pool



Cropping Features: ROI Pool

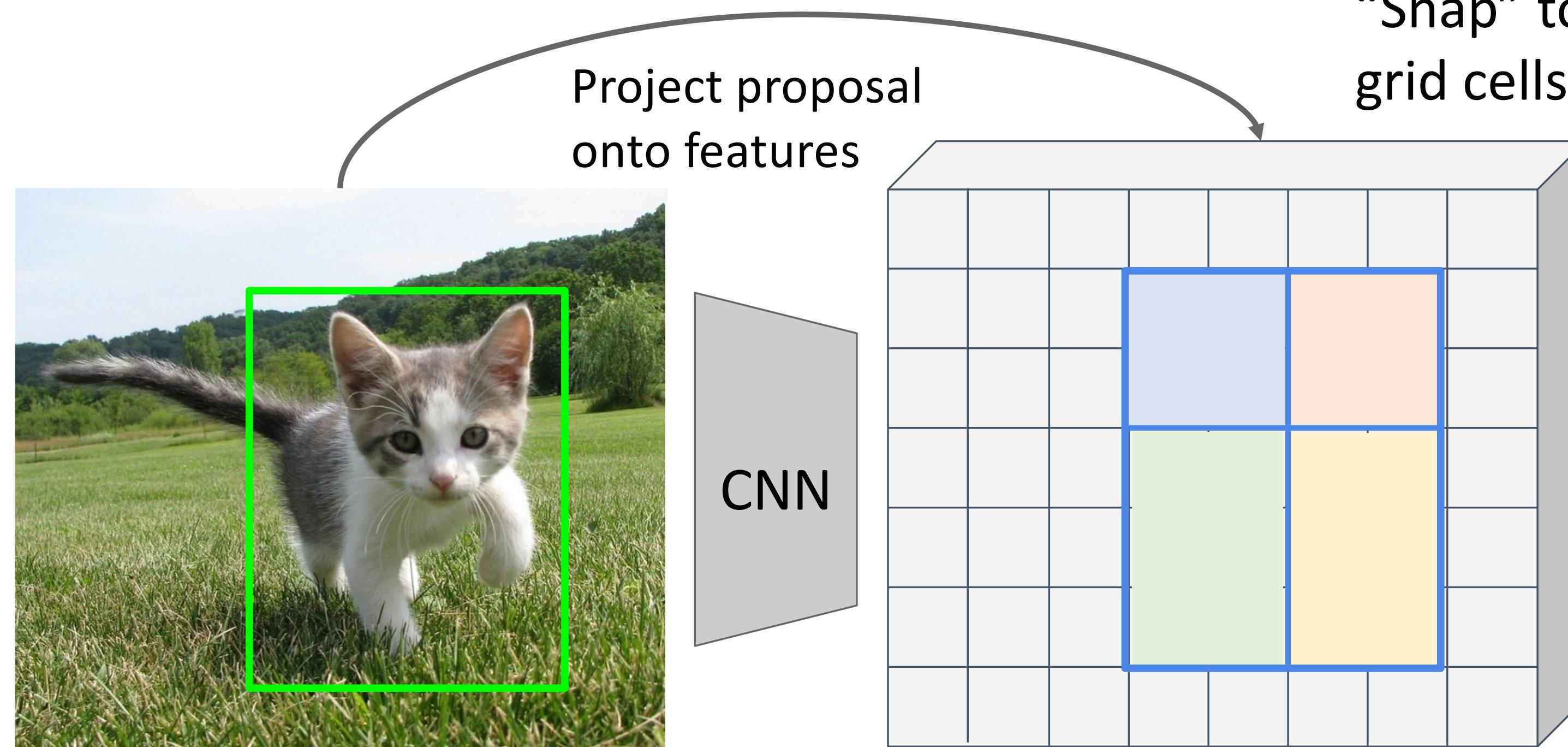


Cropping Features: ROI Pool



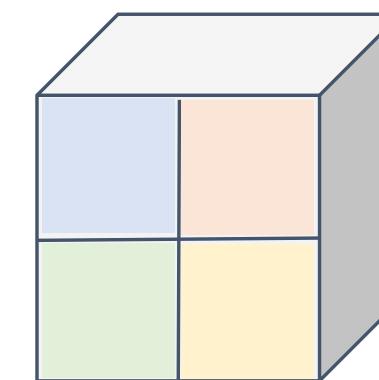
Want features for the box of a fixed size (2x2 in this example, 7x7 or 14x14 in practice)

Cropping Features: ROI Pool



Divide into 2×2 grid of (roughly) equal subregions

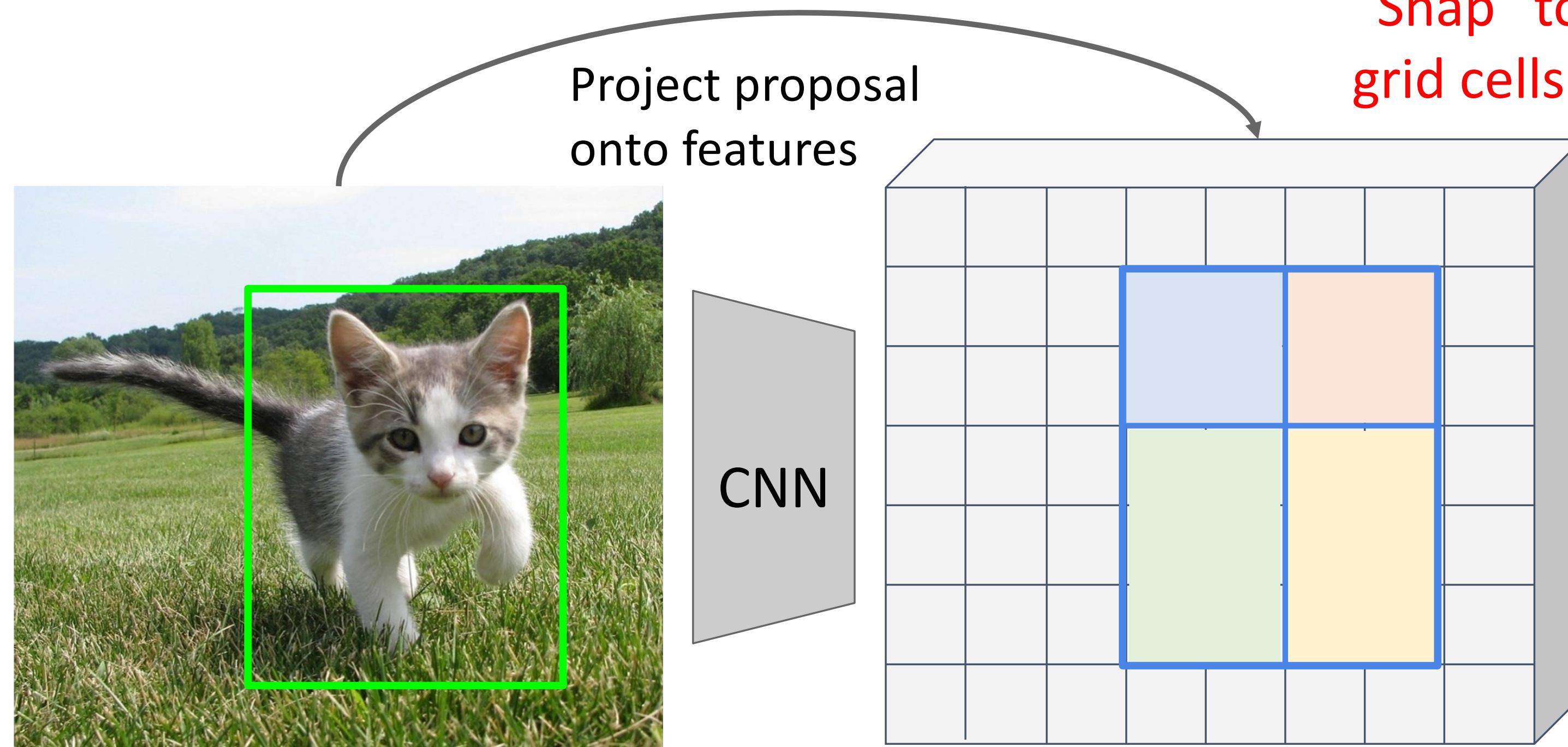
Max-pool within each subregion



Region features
(here $512 \times 2 \times 2$;
In practice $512 \times 7 \times 7$)

Region features always the same size even if input regions have different sizes!

Cropping Features: ROI Pool



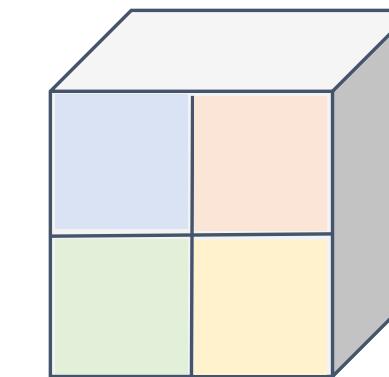
Input Image
(e.g. $3 \times 640 \times 480$)

Image features
(e.g. $512 \times 20 \times 15$)

Problem: Slight misalignment due to snapping; different-sized subregions is weird

Divide into 2×2 grid of (roughly) equal subregions

Max-pool within each subregion

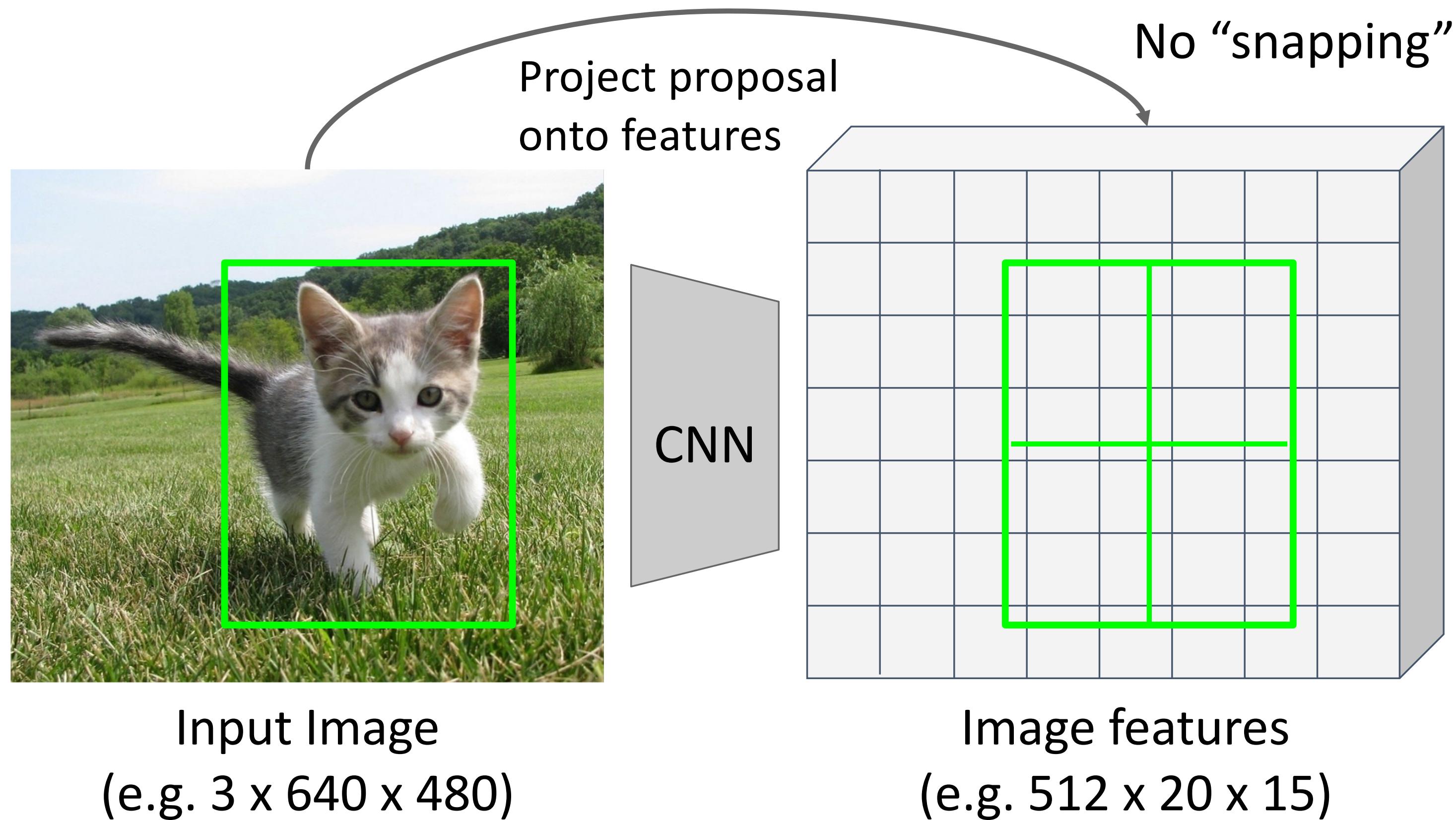


Region features
(here $512 \times 2 \times 2$;
In practice $512 \times 7 \times 7$)

Region features always the same size even if input regions have different sizes!

Cropping Features: ROI Align

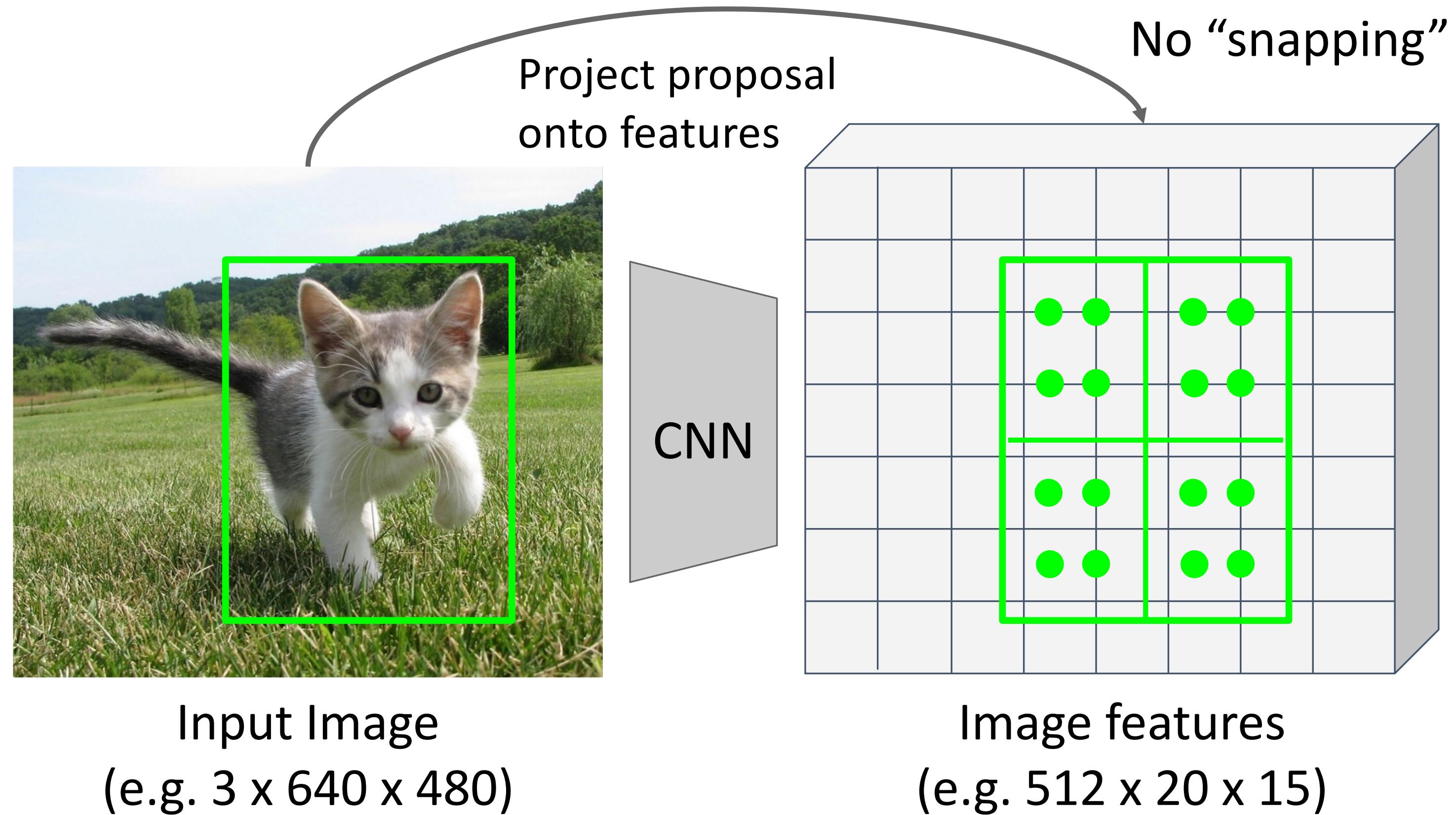
Divide into equal-sized subregions
(may not be aligned to grid!)



Want features for the
box of a fixed size
(2×2 in this example,
 7×7 or 14×14 in practice)

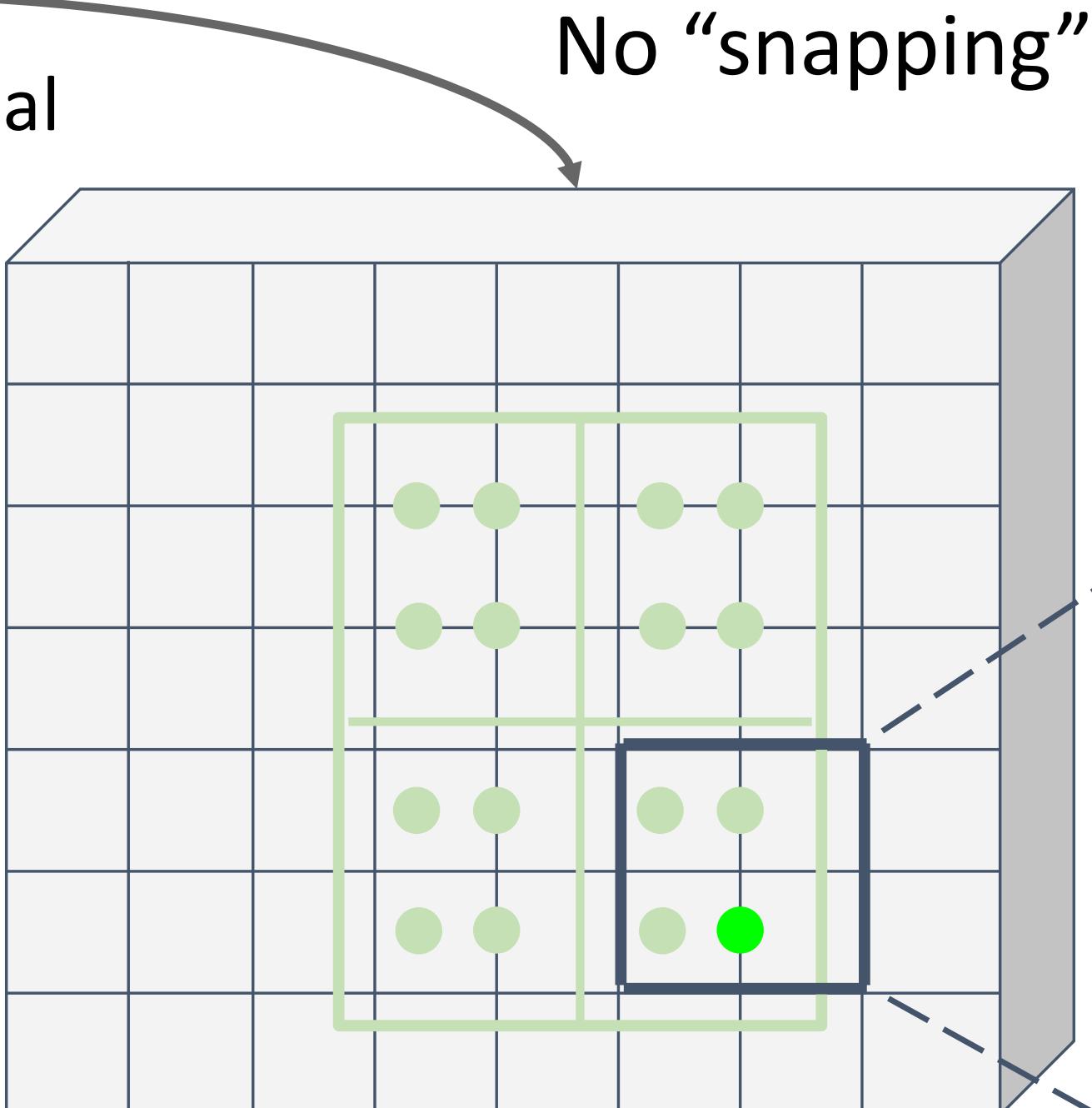
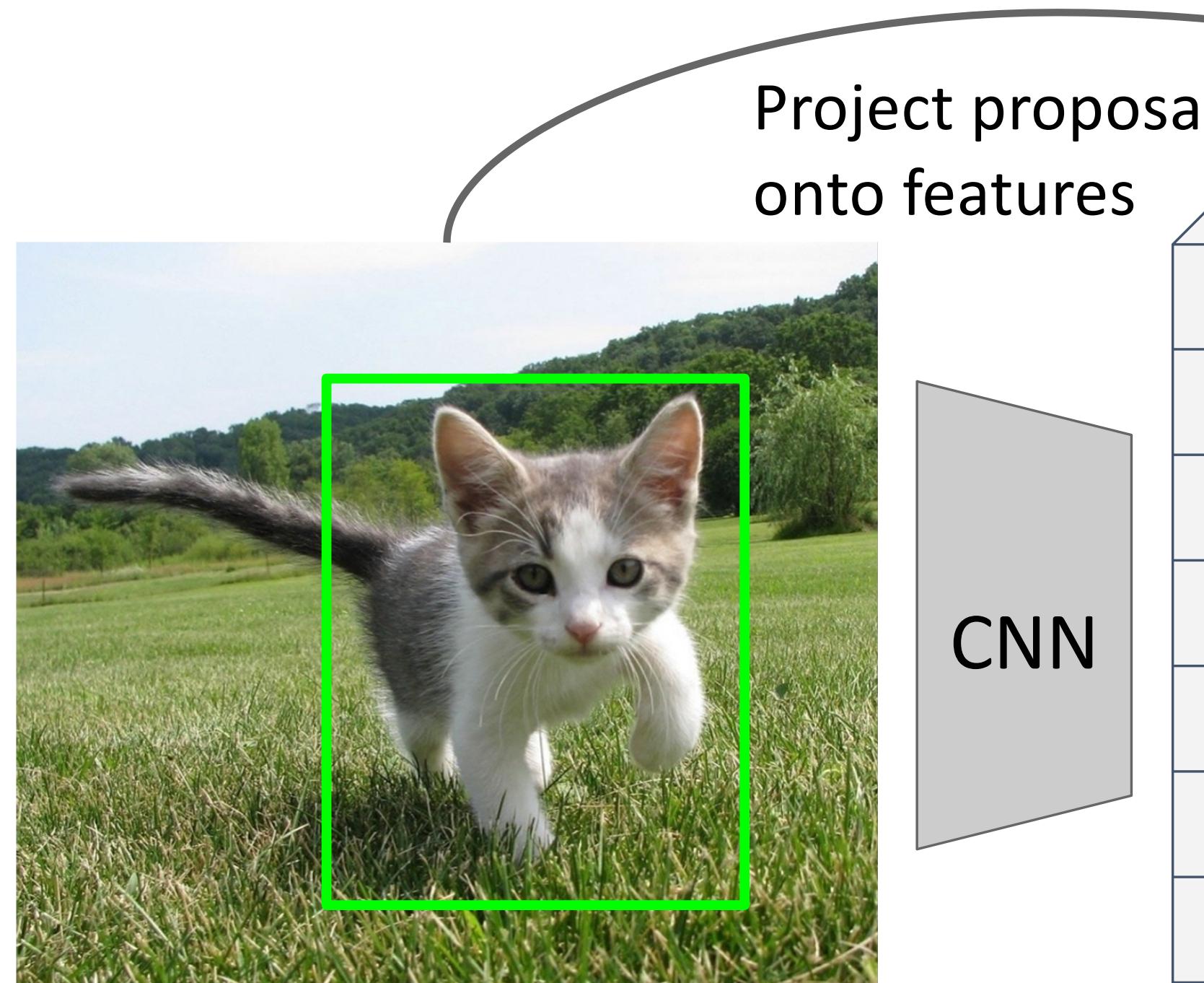
Cropping Features: ROI Align

Divide into equal-sized subregions
(may not be aligned to grid!)

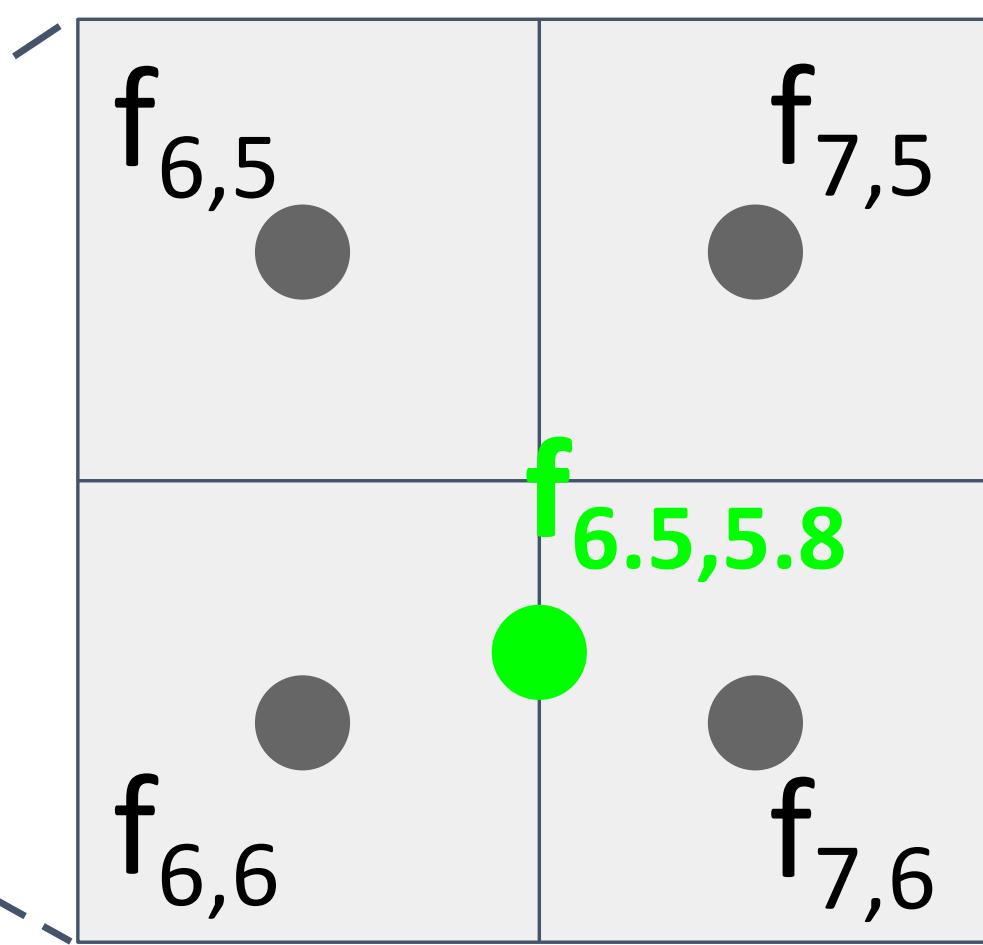


Cropping Features: ROI Align

Divide into equal-sized subregions
(may not be aligned to grid!)



Sample features at regularly-spaced points in each subregion using **bilinear interpolation**

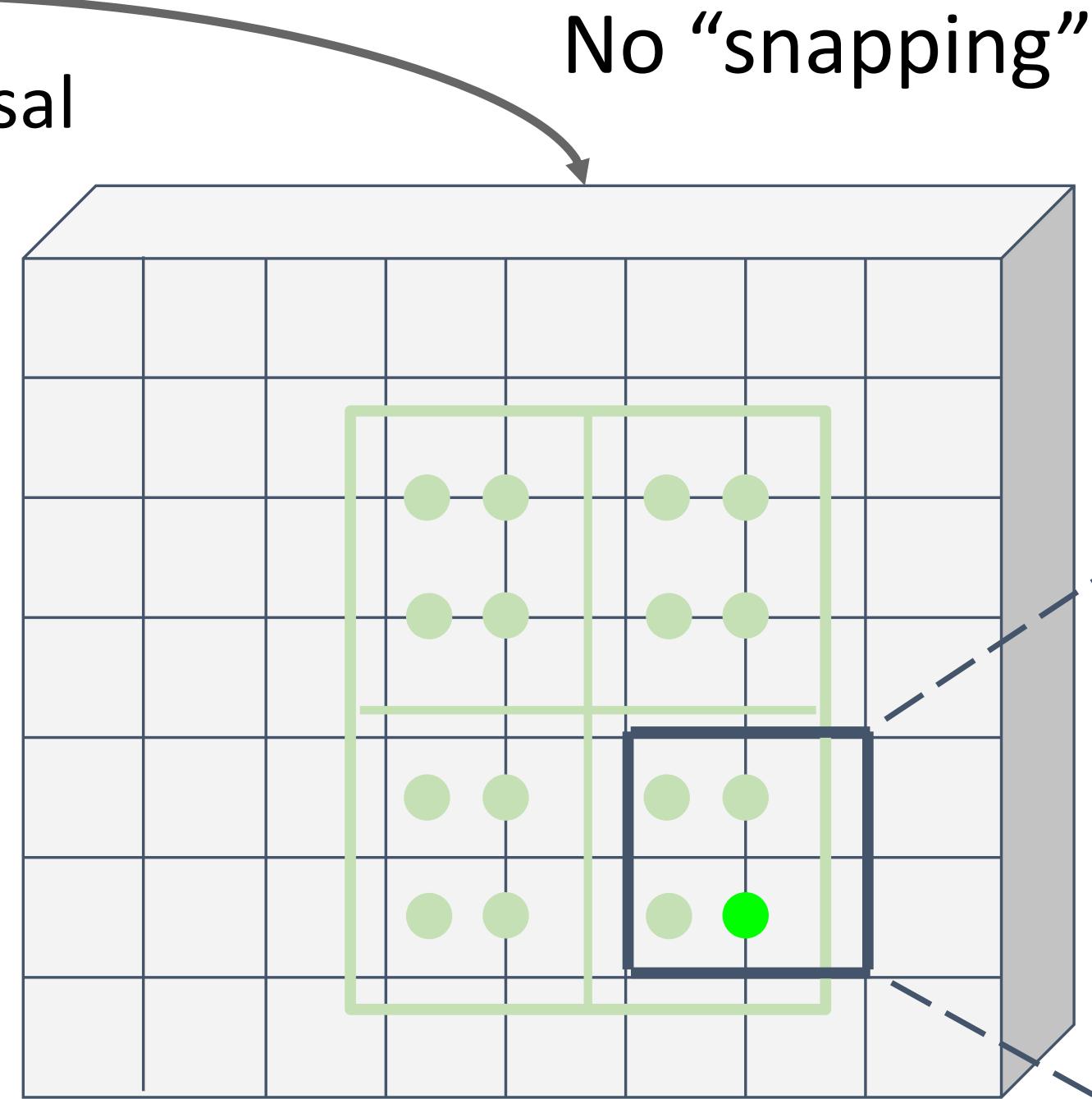
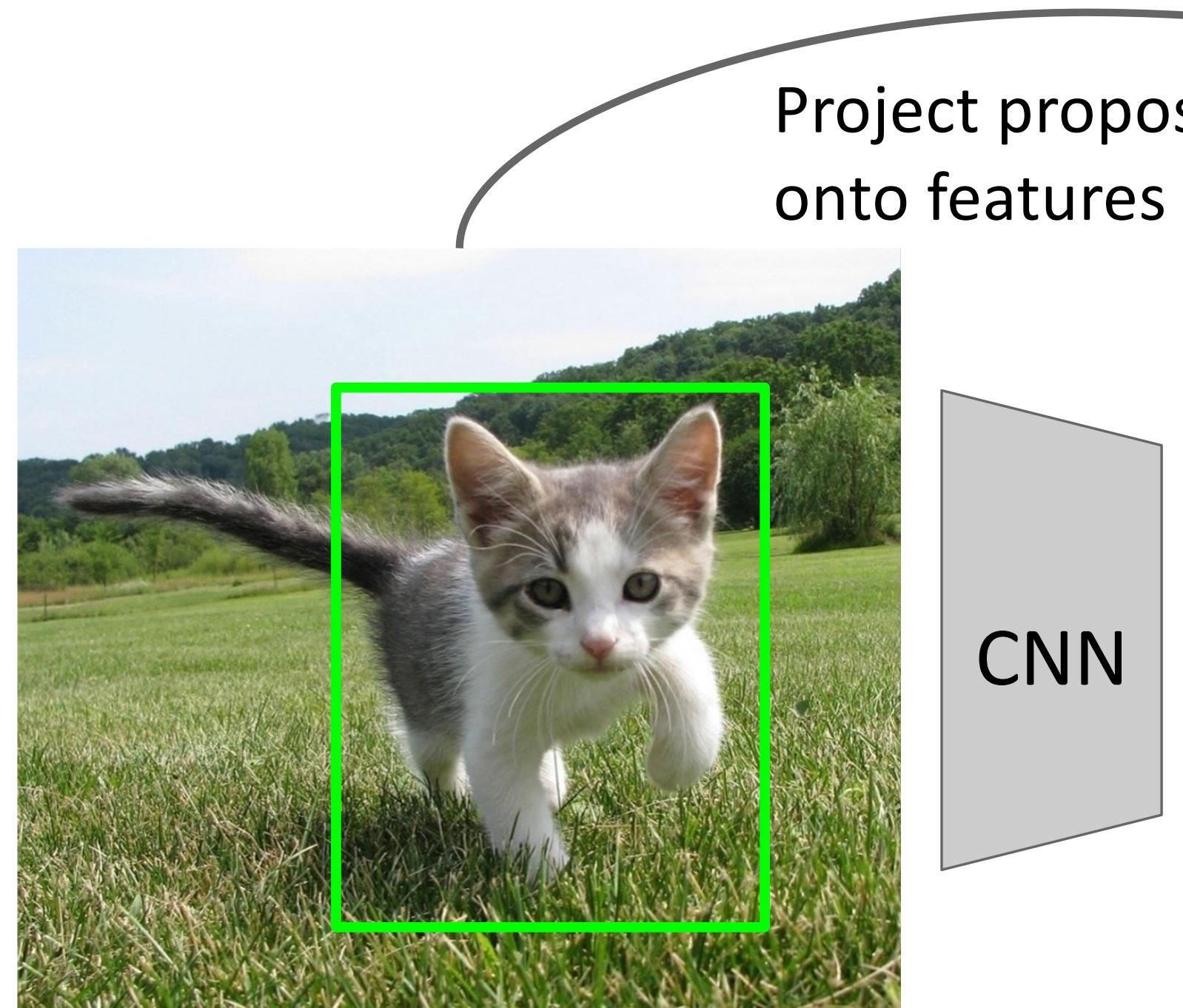


Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

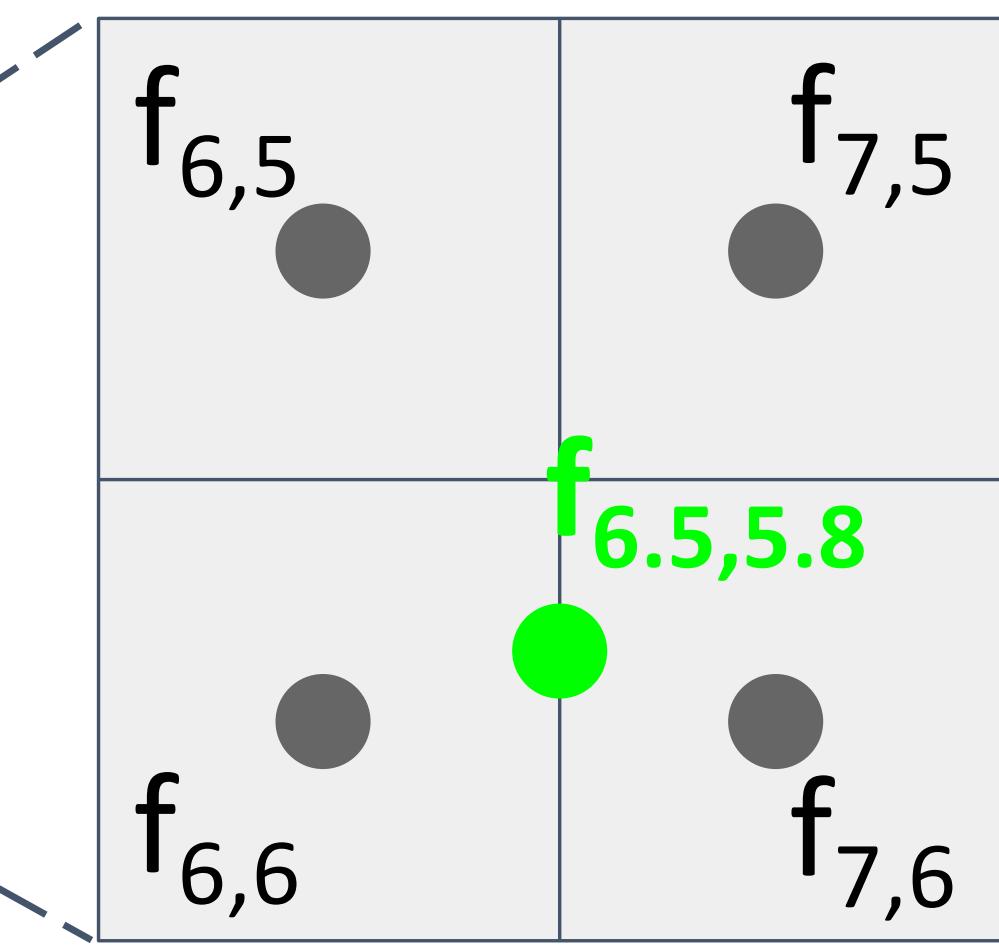
$$f_{xy} = \sum_{i,j=1}^2 f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_j|)$$

Cropping Features: ROI Align

Divide into equal-sized subregions
(may not be aligned to grid!)



Sample features at regularly-spaced points in each subregion using **bilinear interpolation**

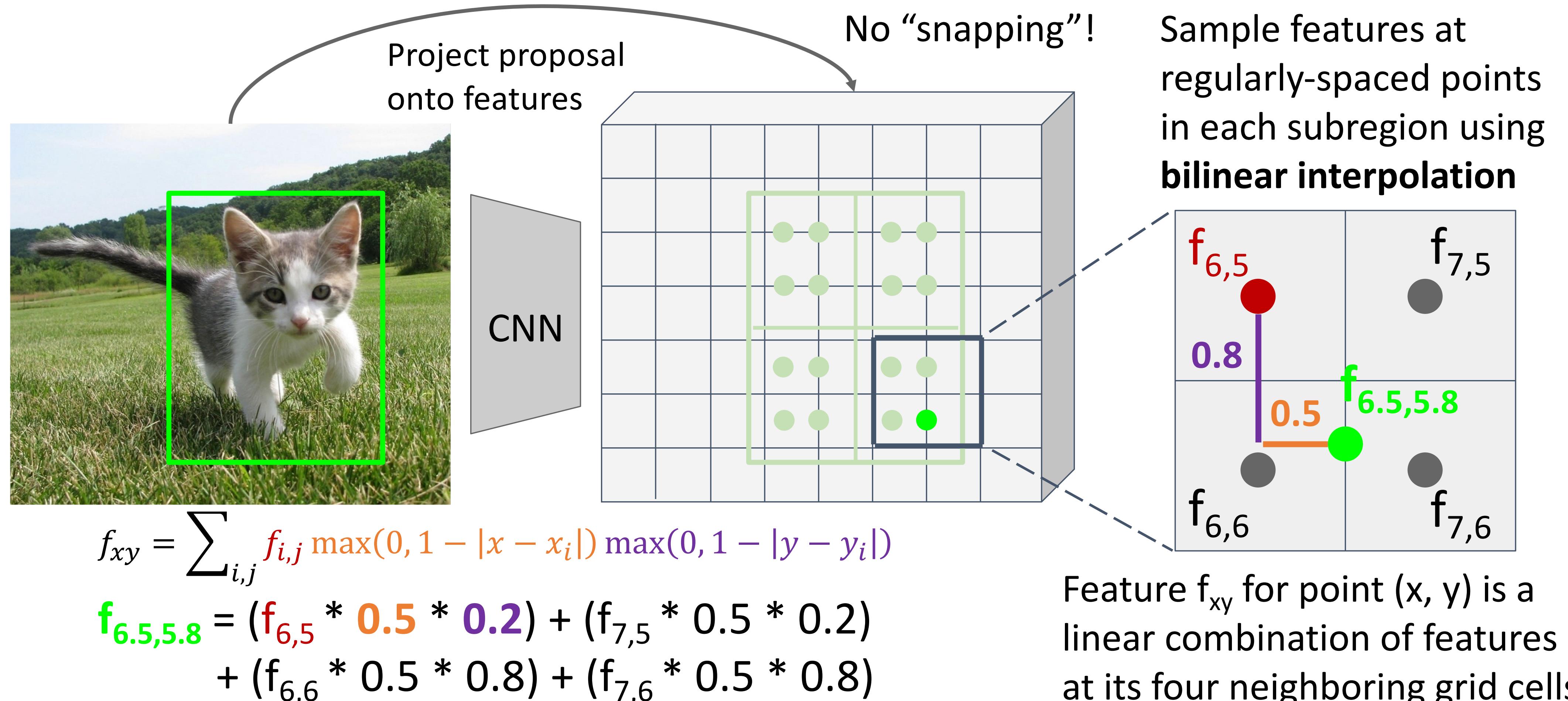


Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

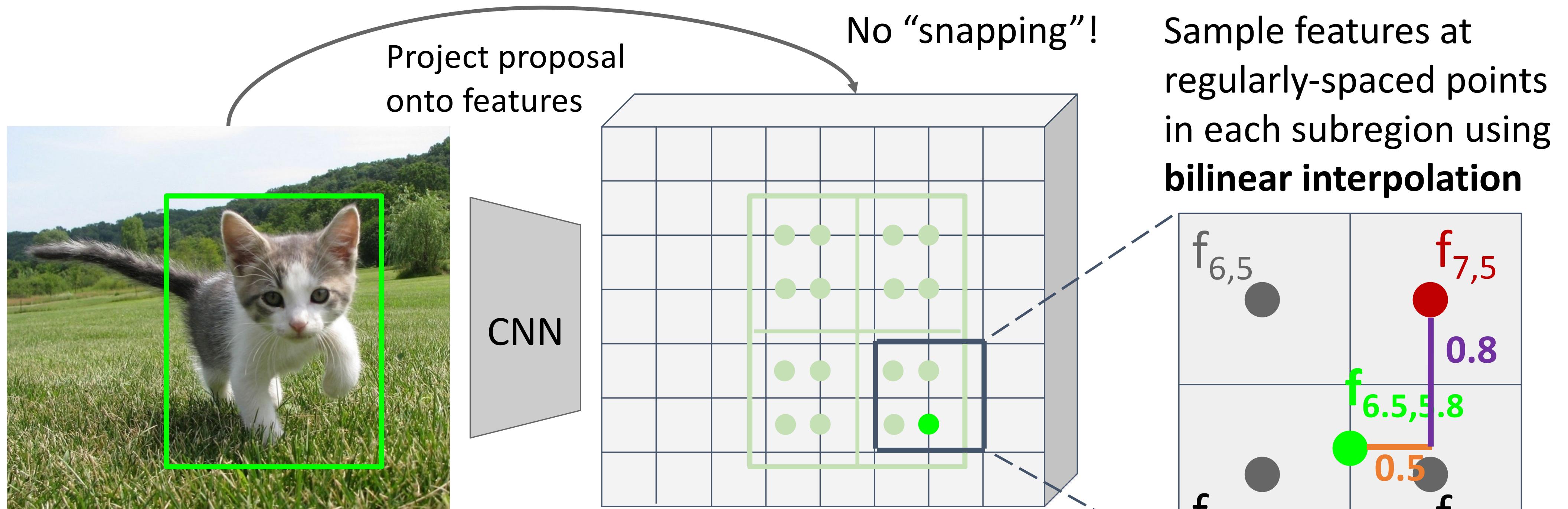
$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

$$\begin{aligned} f_{6.5,5.8} &= (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) \\ &\quad + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8) \end{aligned}$$

Cropping Features: ROI Align

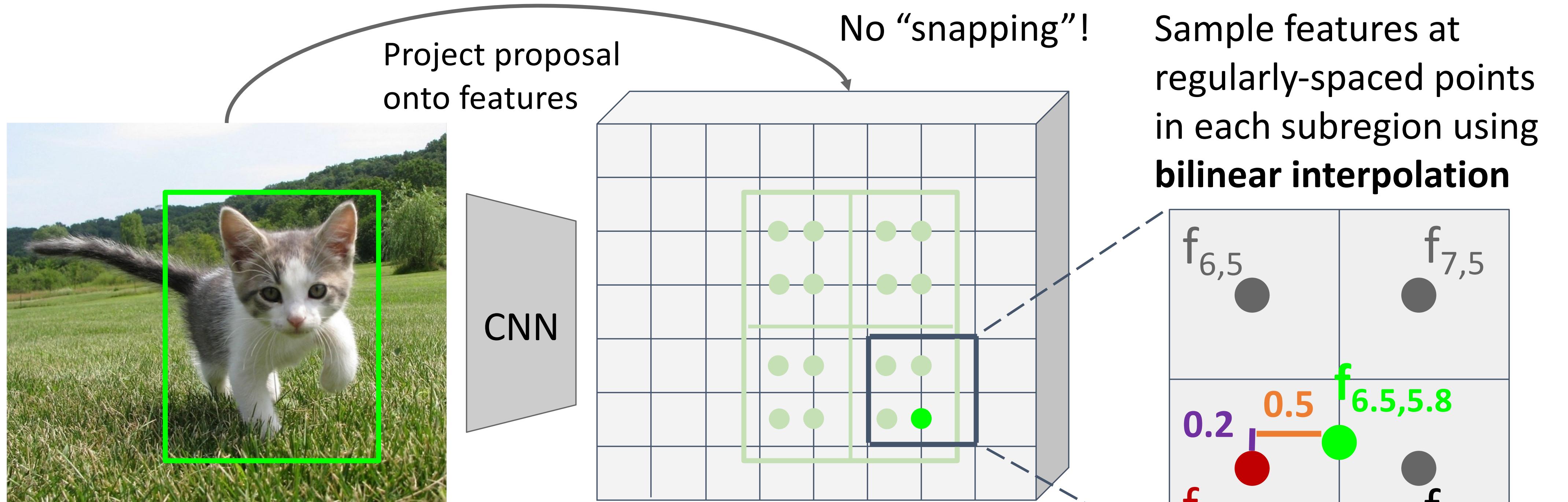


Cropping Features: ROI Align



Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

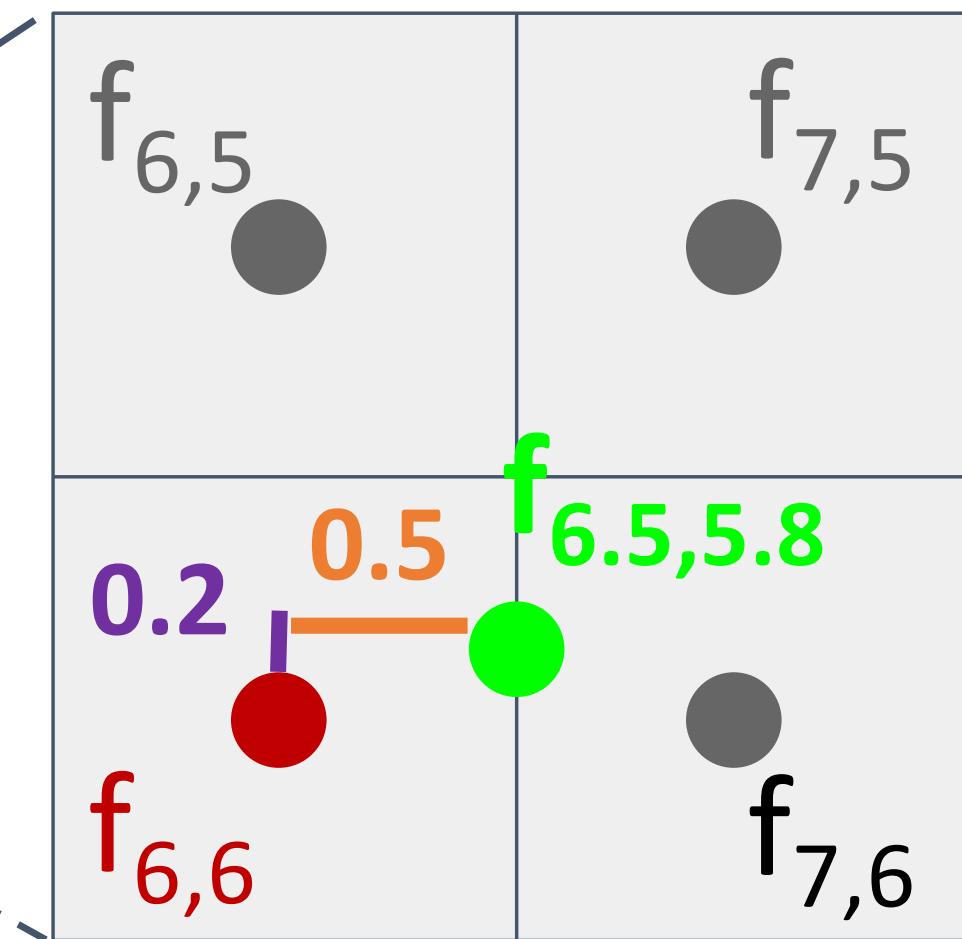
Cropping Features: ROI Align



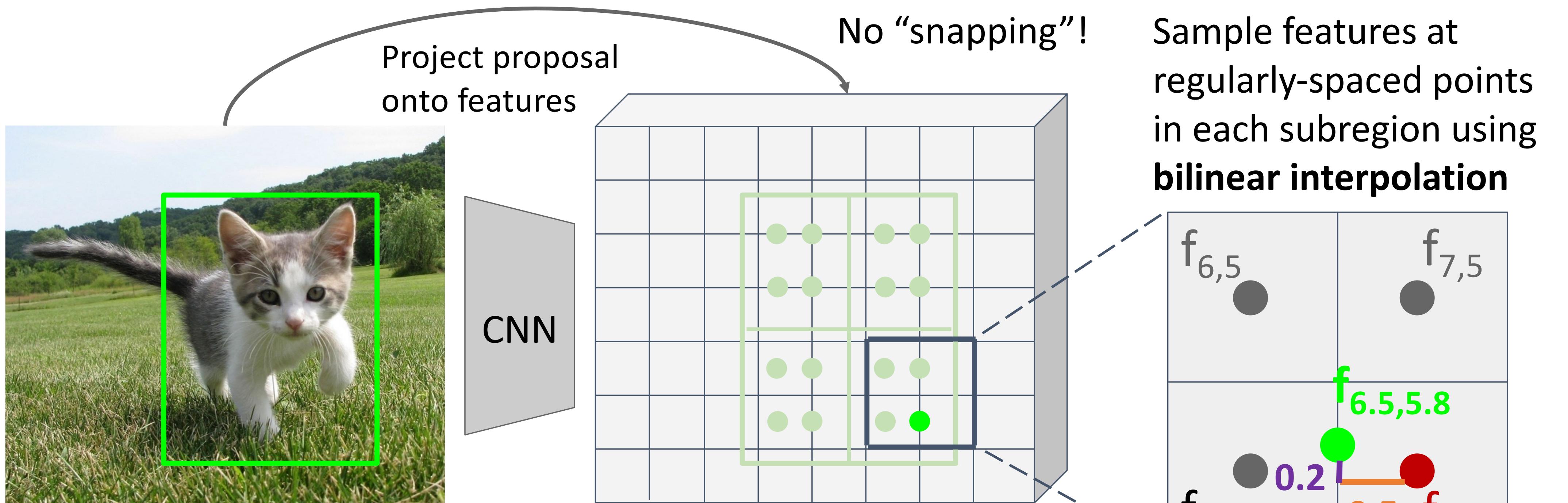
$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

$$\begin{aligned} f_{6.5,5.8} &= (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) \\ &\quad + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8) \end{aligned}$$

Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:



Cropping Features: ROI Align

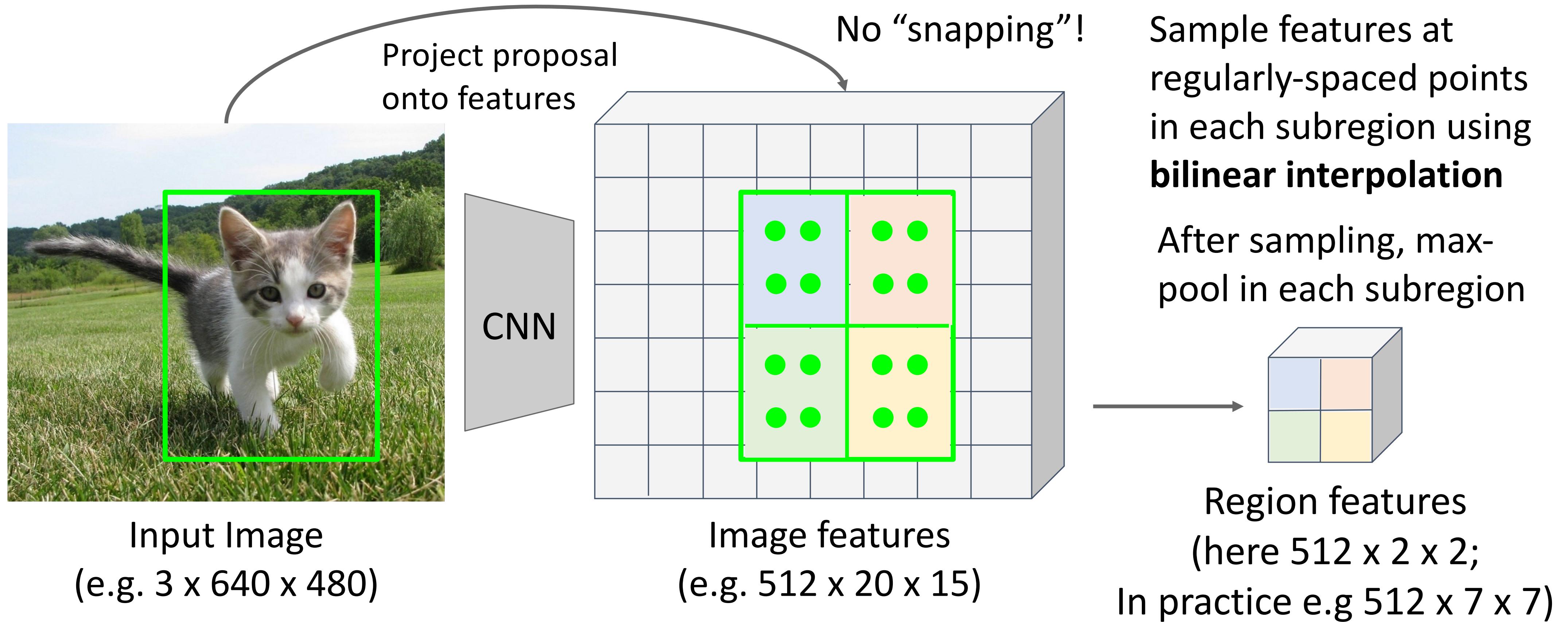


$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

$$\begin{aligned} f_{6,5,8} &= (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) \\ &\quad + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8) \end{aligned}$$

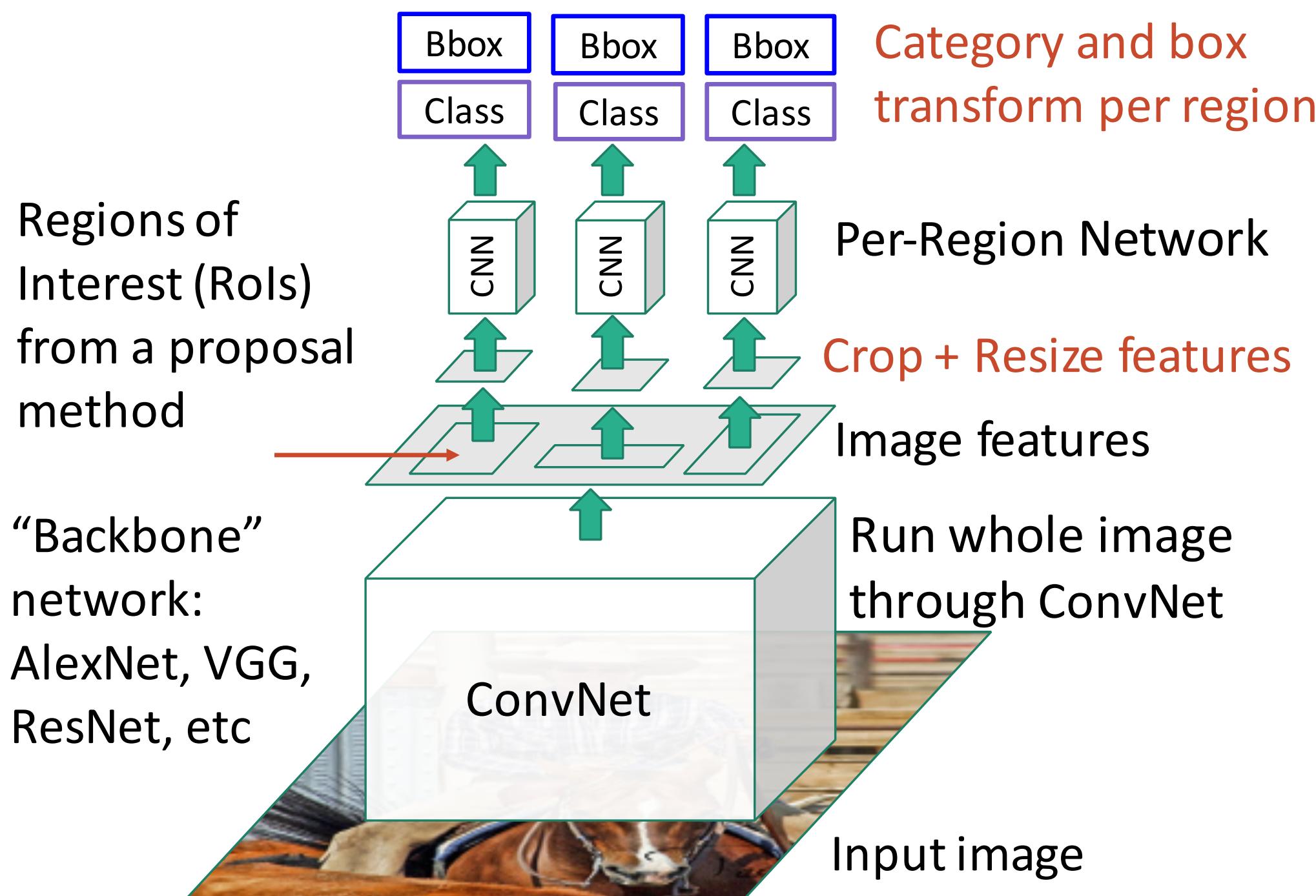
Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Cropping Features: ROI Align

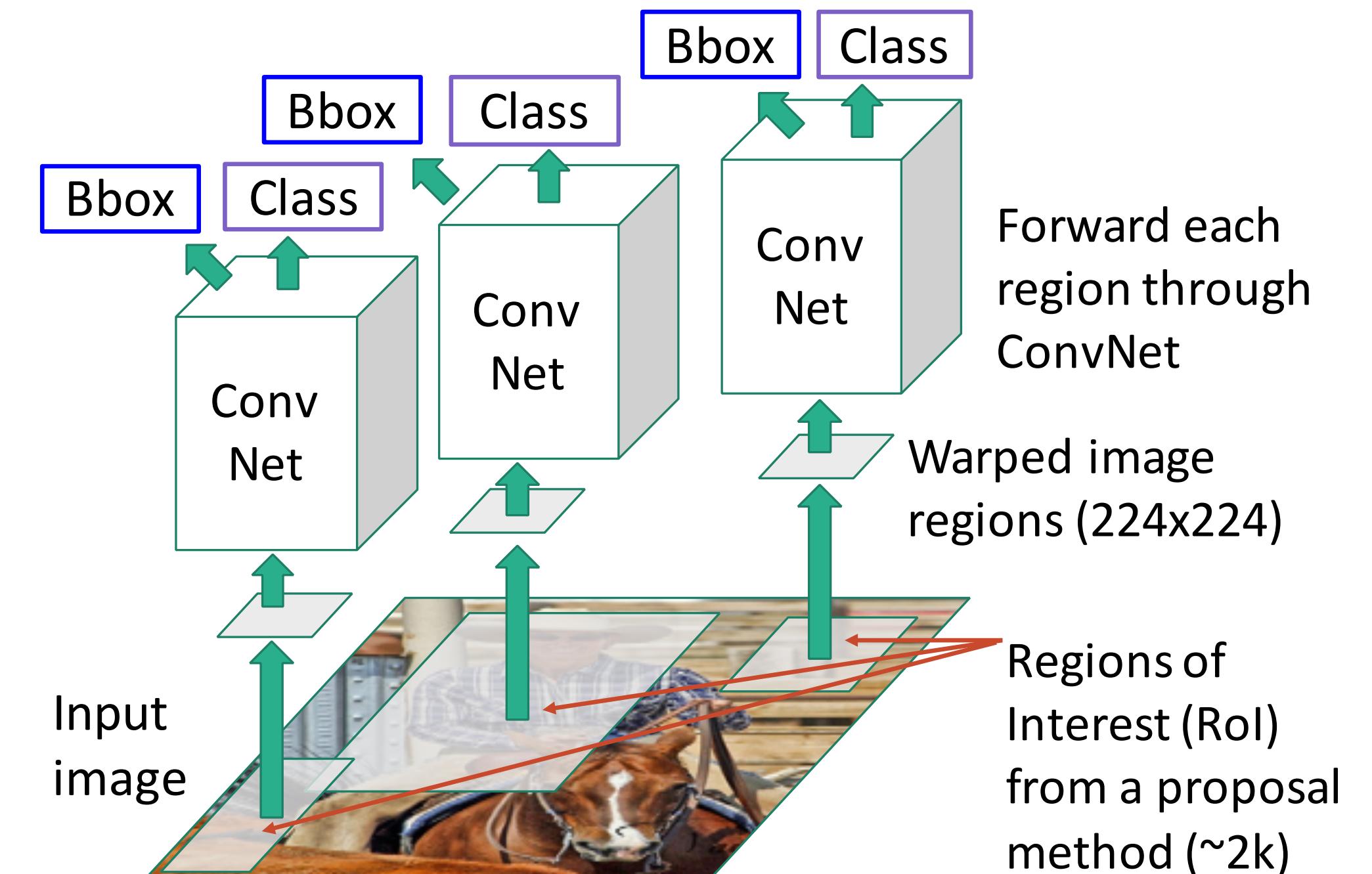


Fast R-CNN vs “Slow” R-CNN

Fast R-CNN: Apply differentiable cropping to shared image features

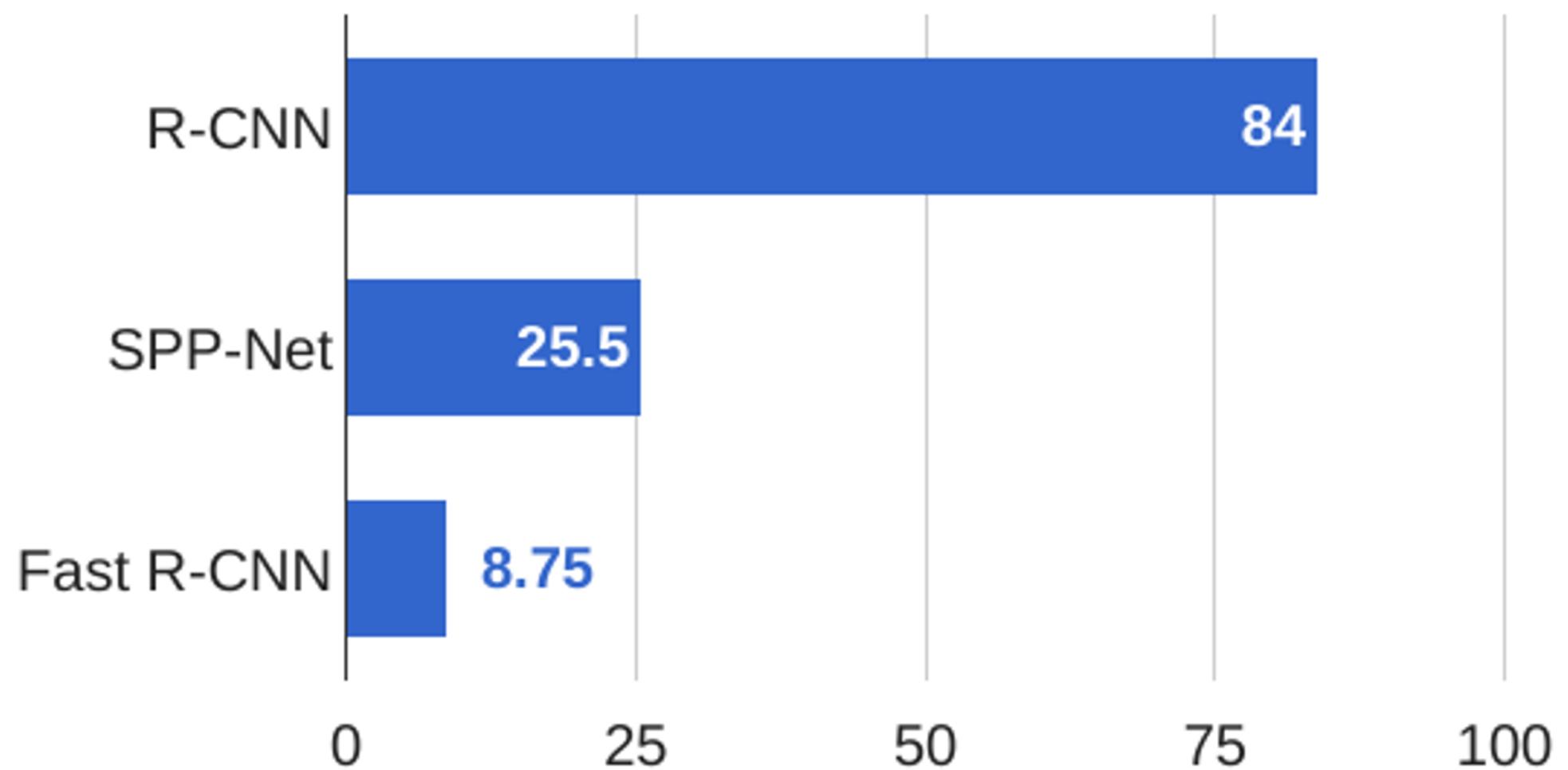


“Slow” R-CNN: Apply differentiable cropping to shared image features

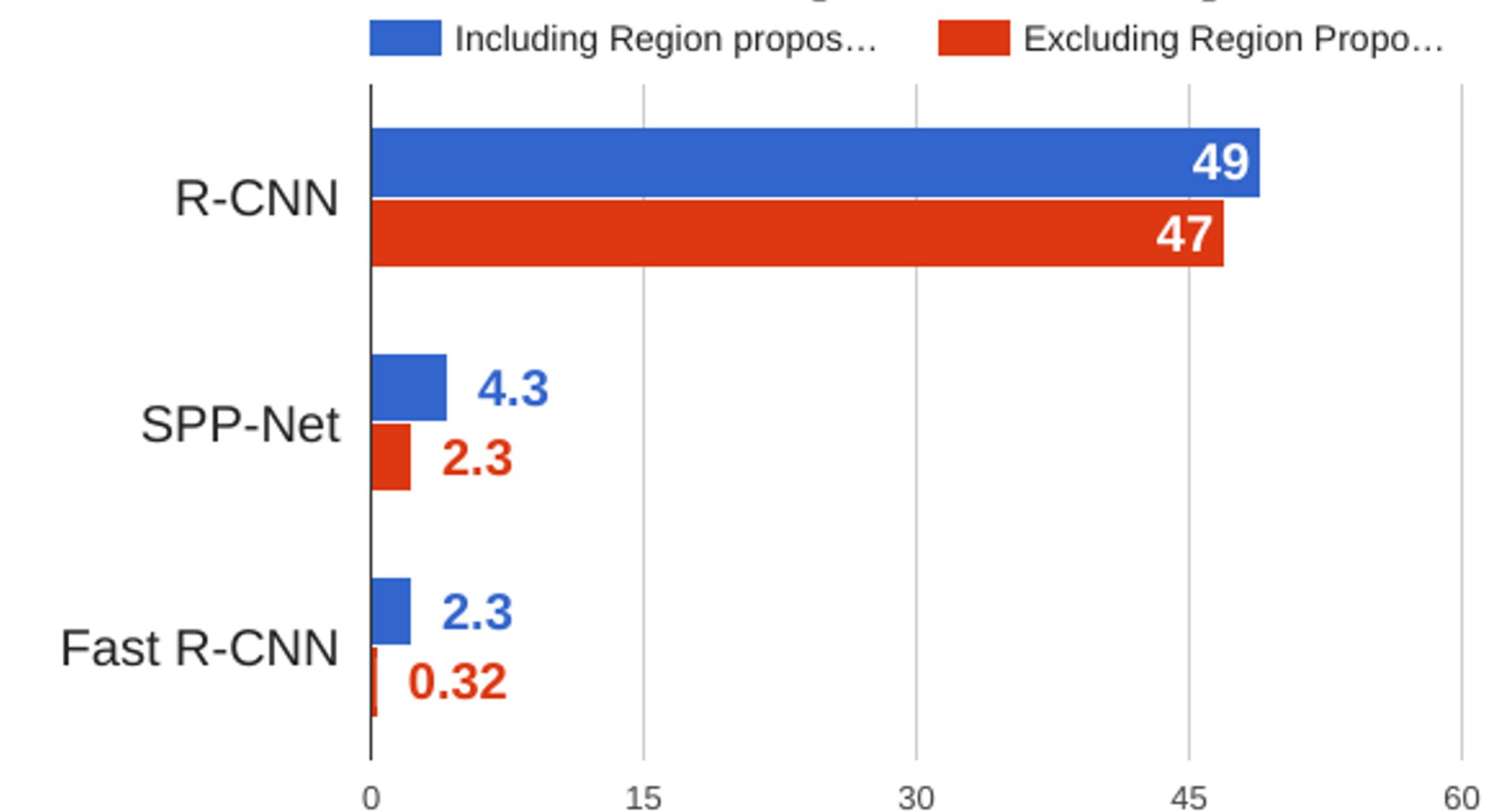


Fast R-CNN vs “Slow” R-CNN

Training time (Hours)



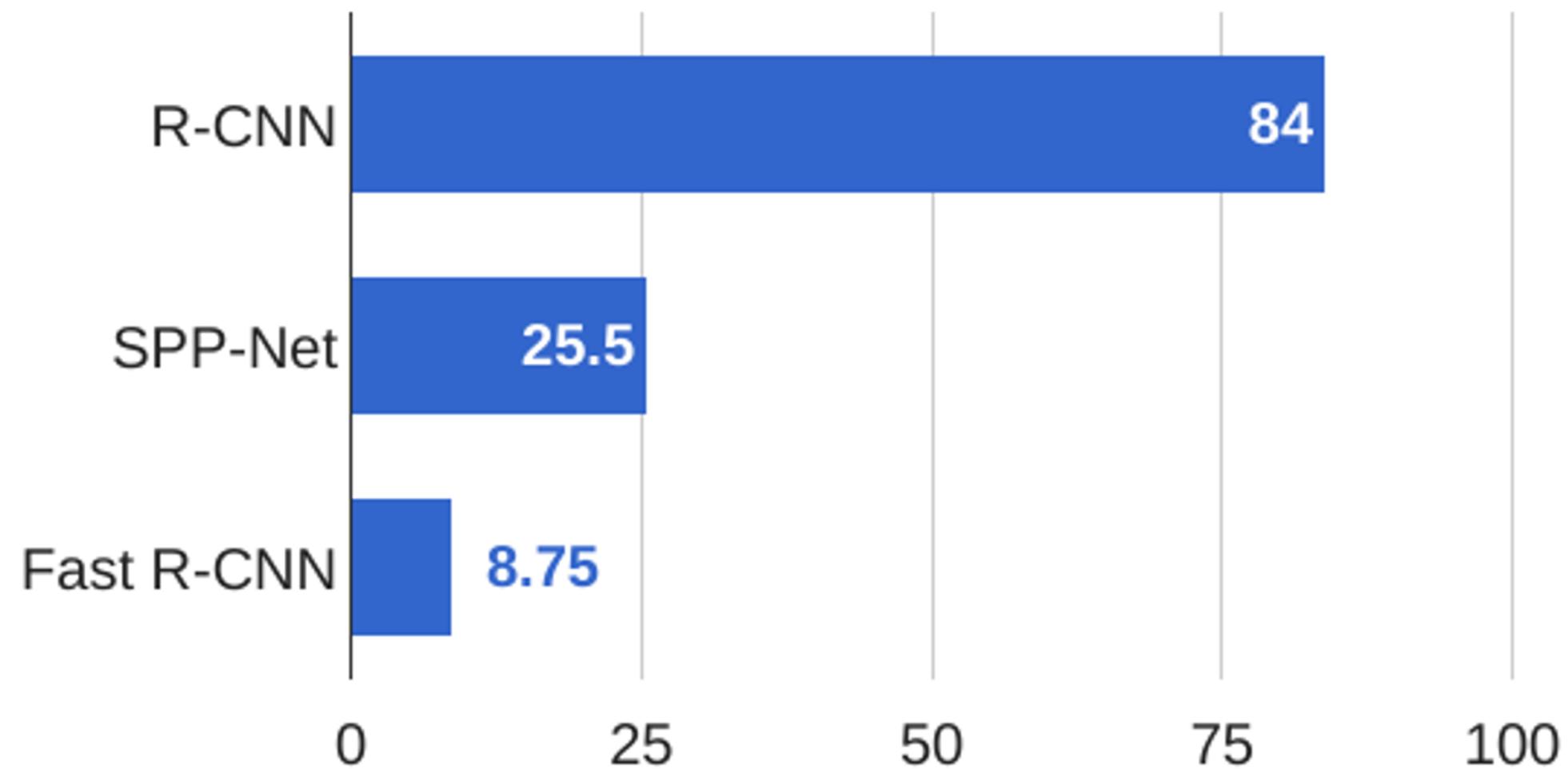
Test time (seconds)



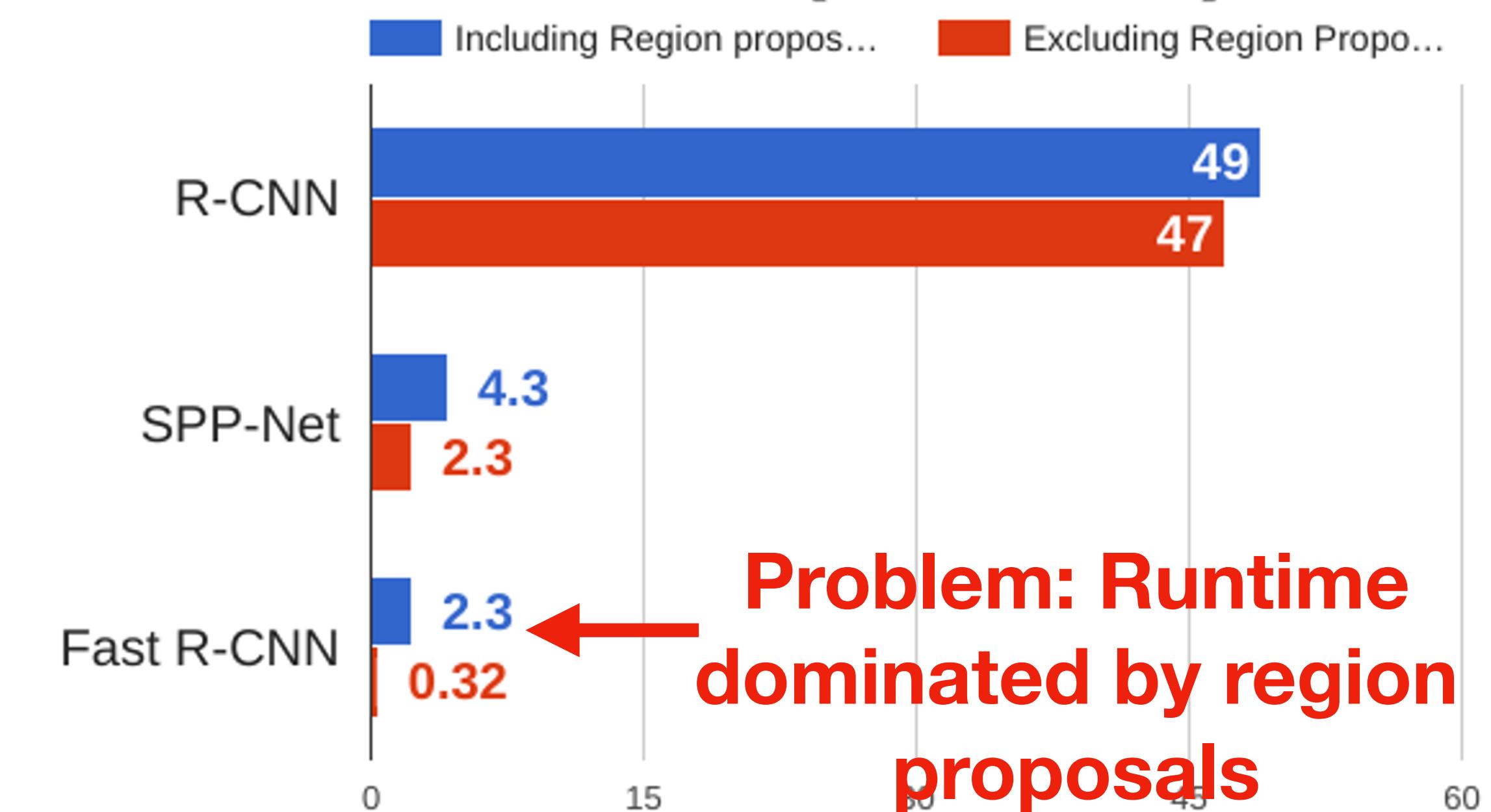
Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.
He et al, “Spatial pyramid pooling in deep convolutional networks for visual recognition”, ECCV 2014
Girshick, “Fast R-CNN”, ICCV 2015

Fast R-CNN vs “Slow” R-CNN

Training time (Hours)

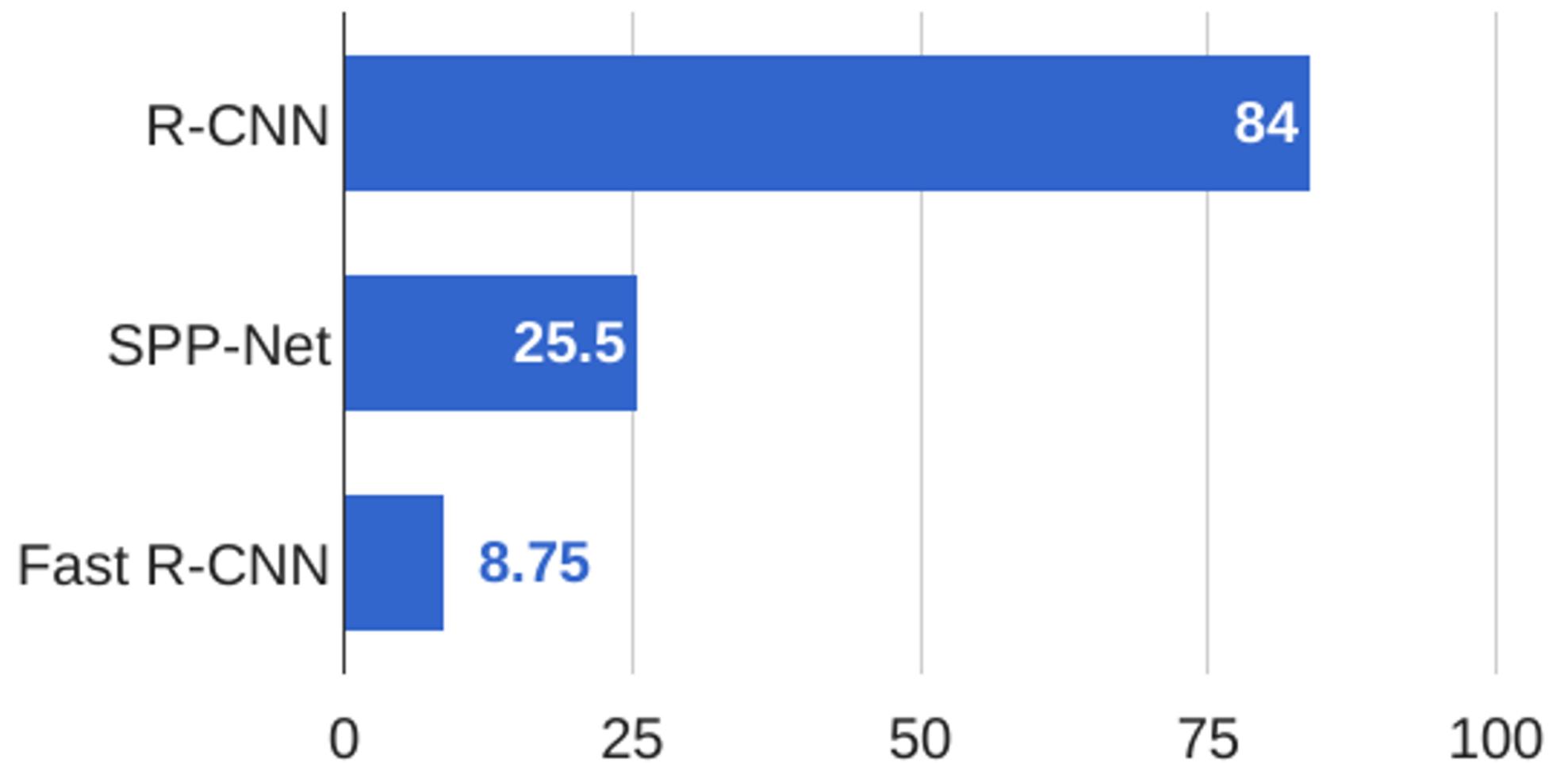


Test time (seconds)

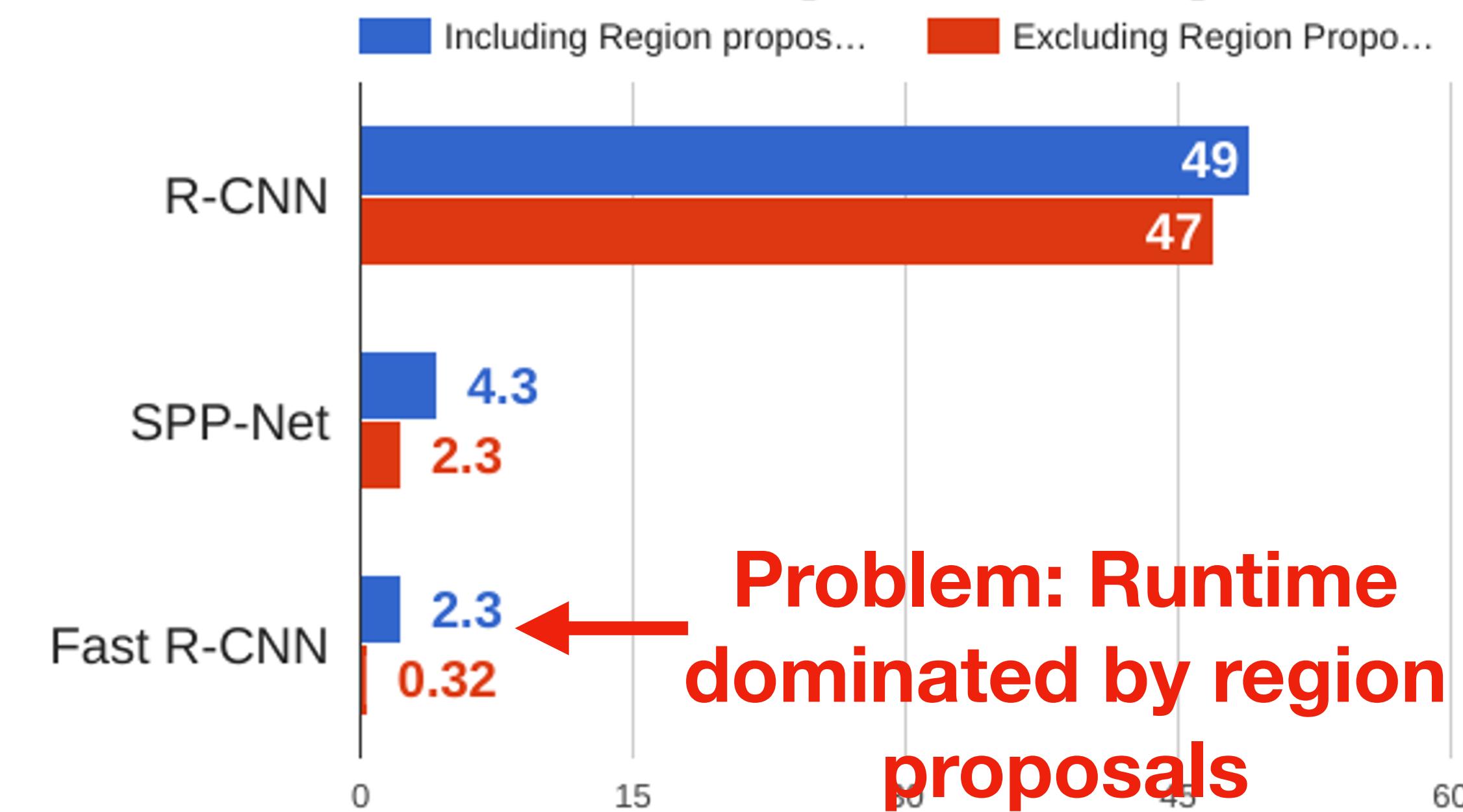


Fast R-CNN vs “Slow” R-CNN

Training time (Hours)



Test time (seconds)

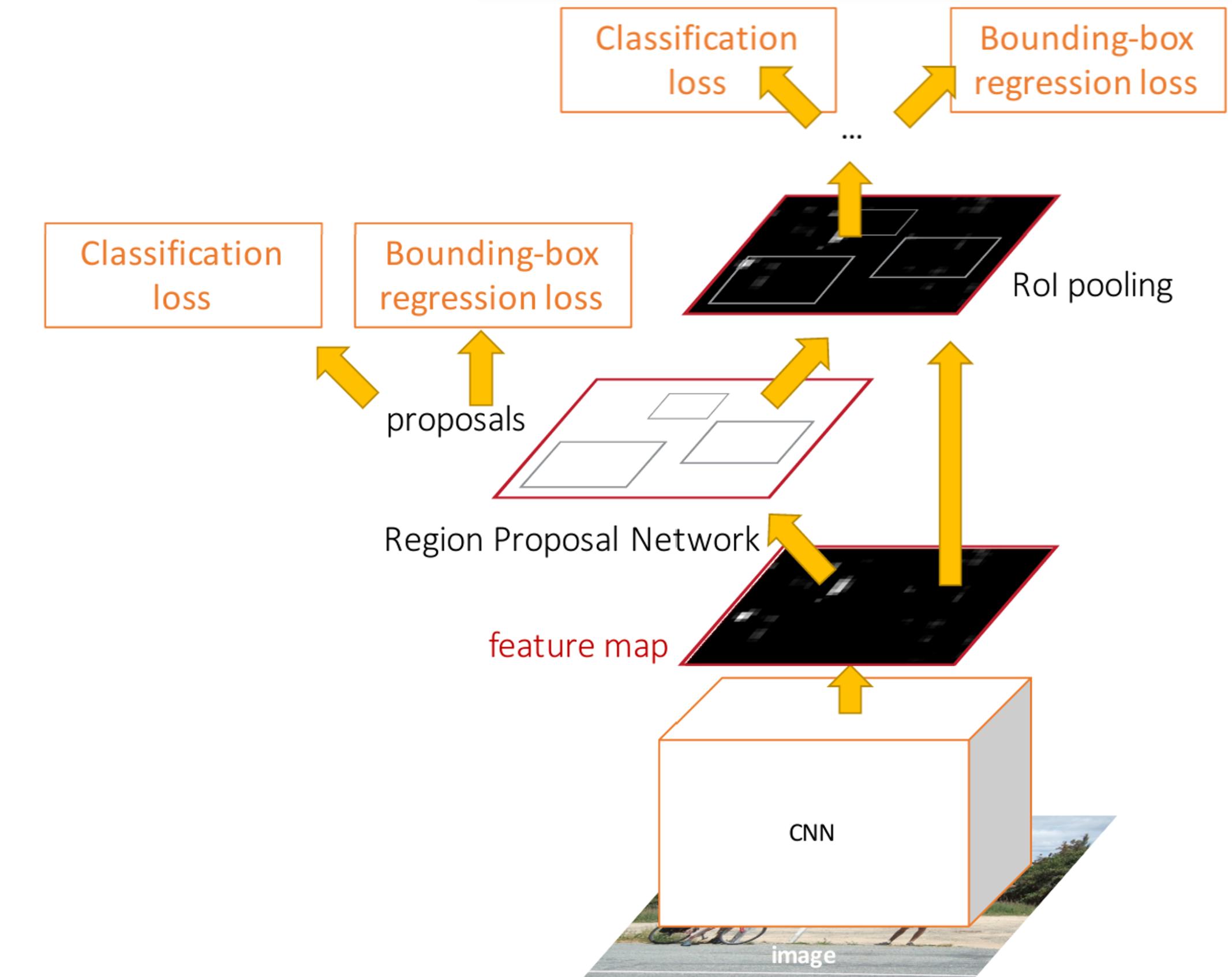


Recall: Region proposals computed by heuristic “Selective search” algorithm on CPU – let’s learn them with a CNN

Faster R-CNN: Learnable Region Proposals

Insert **Region Proposal Network (RPN)** to predict proposals from features

Otherwise same as Fast R-CNN:
Crop features for each proposal,
classify each one



Region Proposal Network (RPN)

Run backbone CNN to get
features aligned to input image



Input Image
(e.g. $3 \times 640 \times 480$)

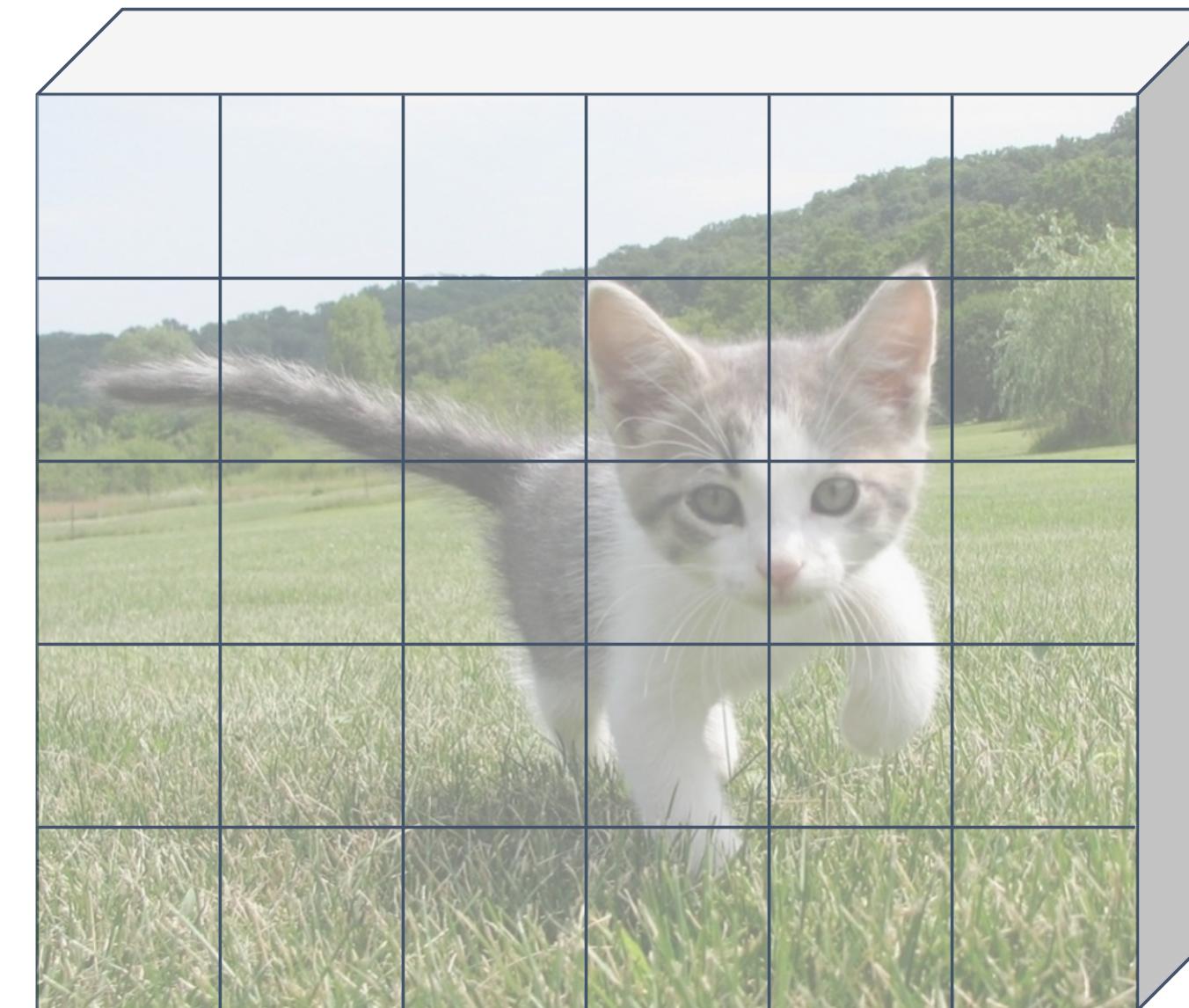
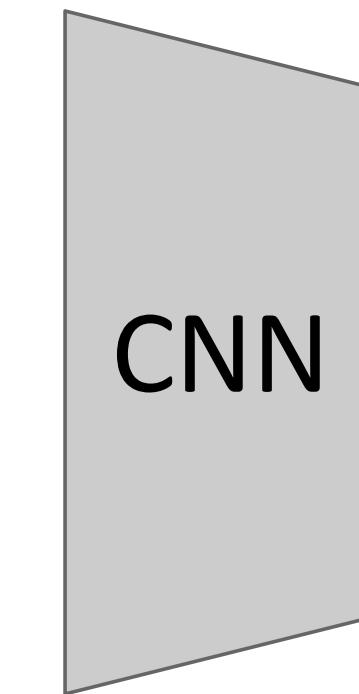


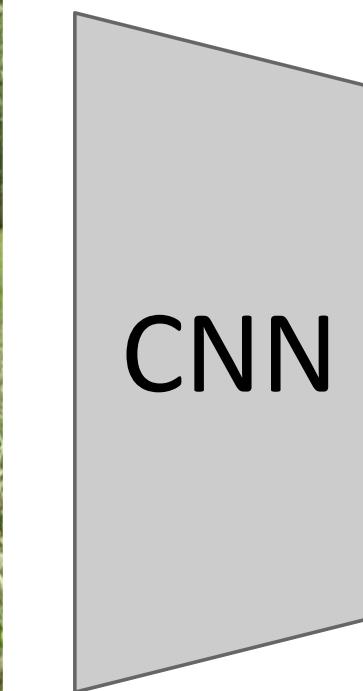
Image features
(e.g. $512 \times 5 \times 6$)

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. $3 \times 640 \times 480$)



Each feature corresponds to a point in the input

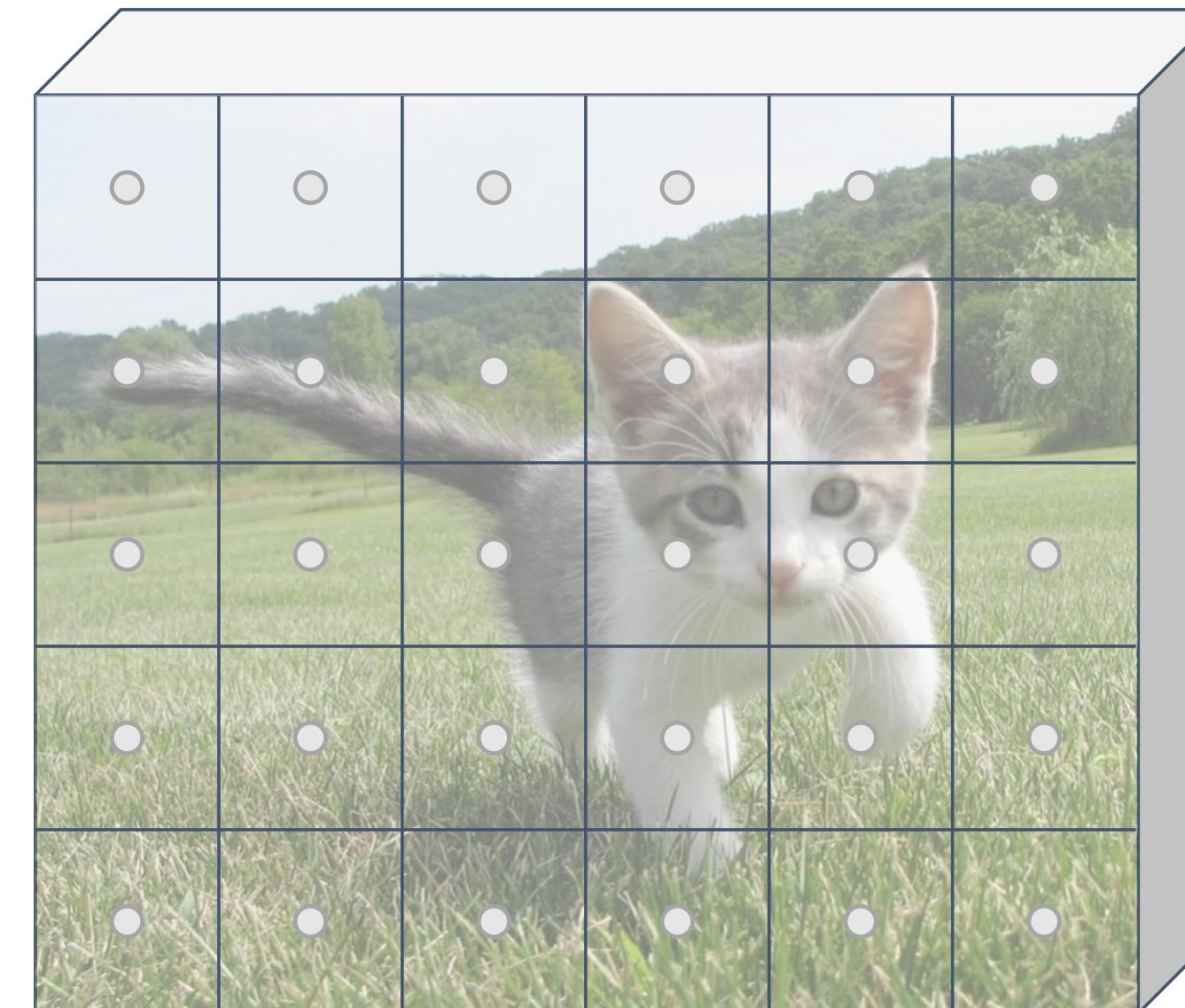
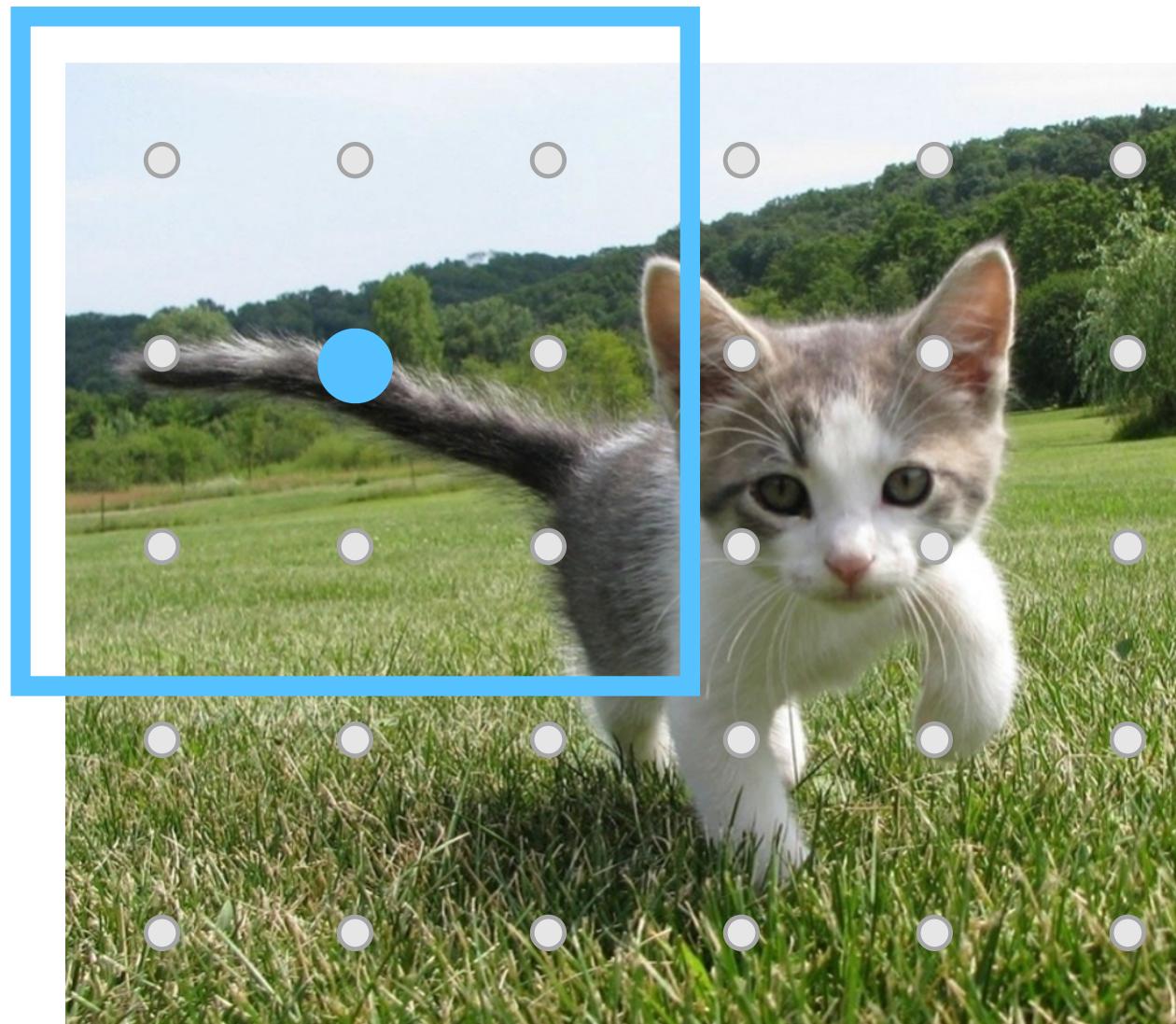


Image features
(e.g. $512 \times 5 \times 6$)

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. $3 \times 640 \times 480$)

Each feature corresponds to a point in the input

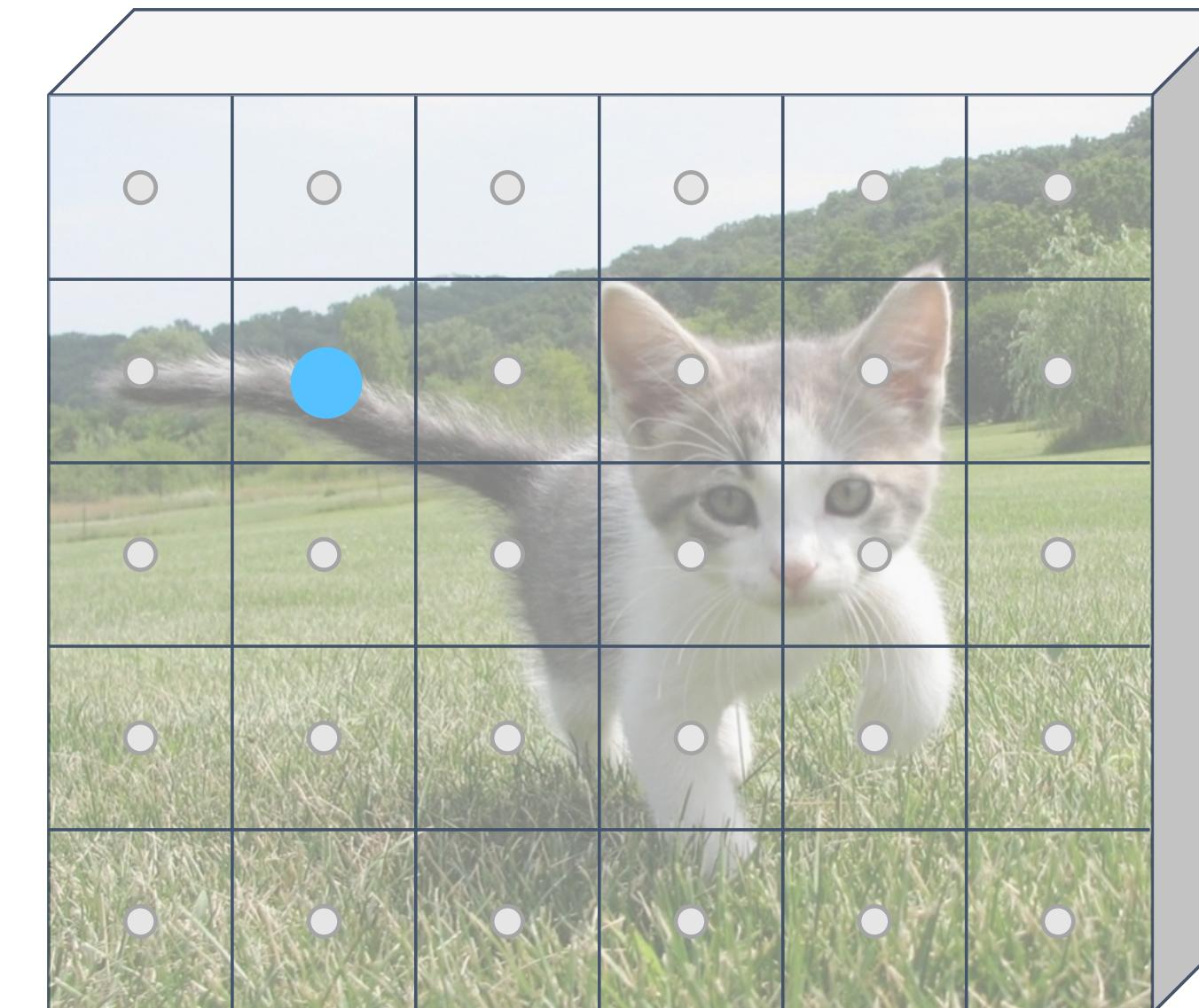
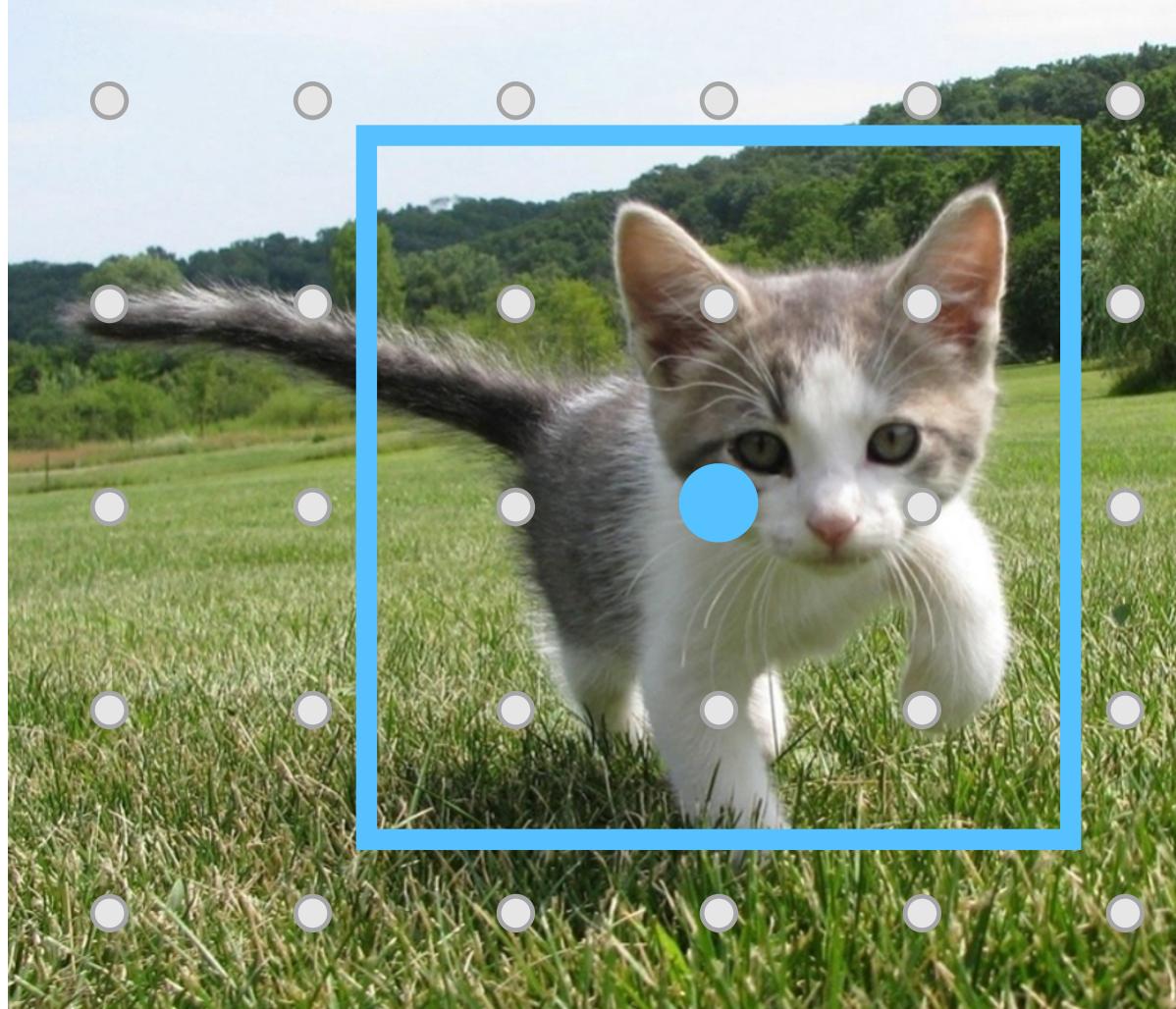


Image features
(e.g. $512 \times 5 \times 6$)

Imagine an **anchor box** of fixed size at each point in the feature map

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. $3 \times 640 \times 480$)

Each feature corresponds to a point in the input

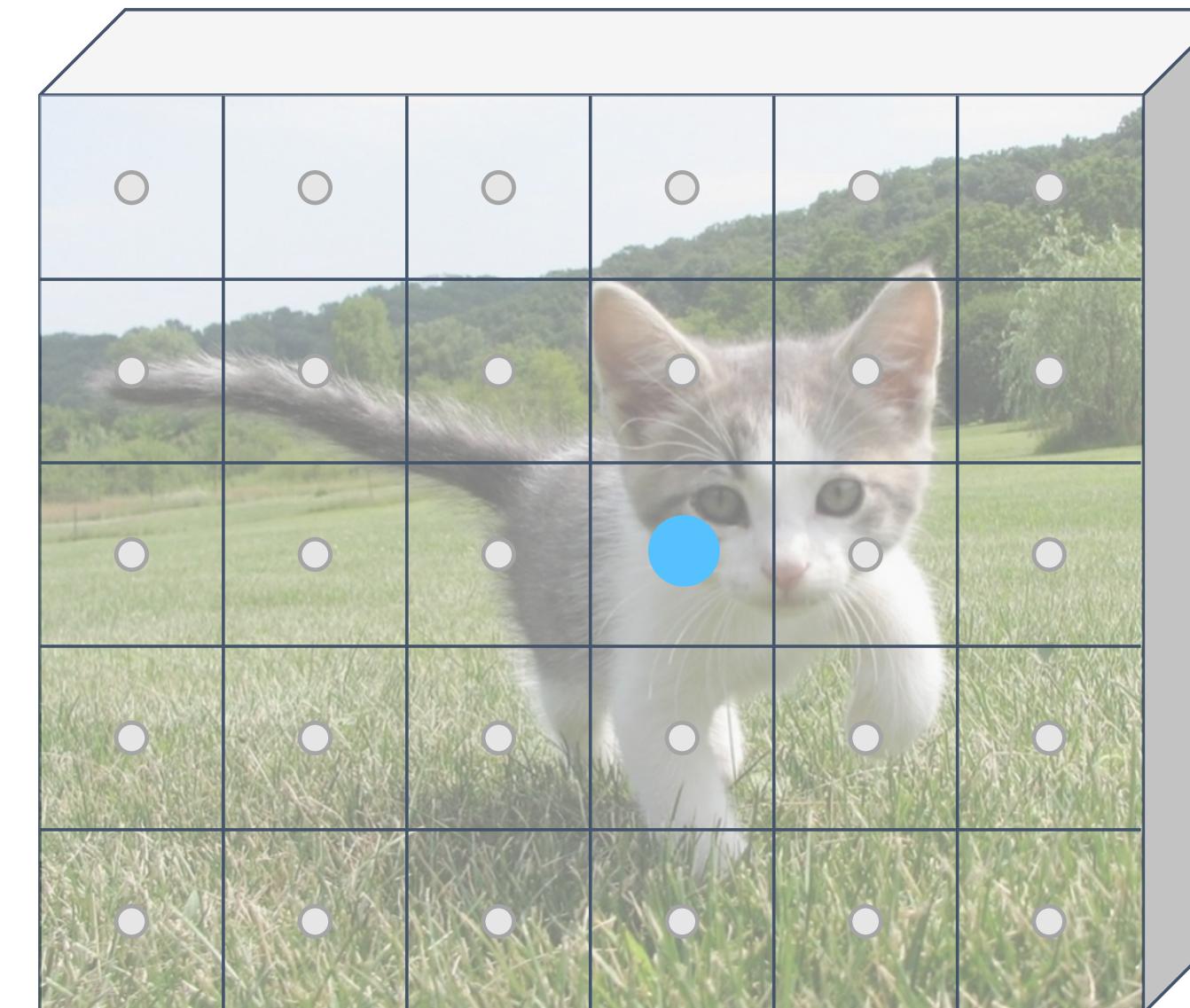
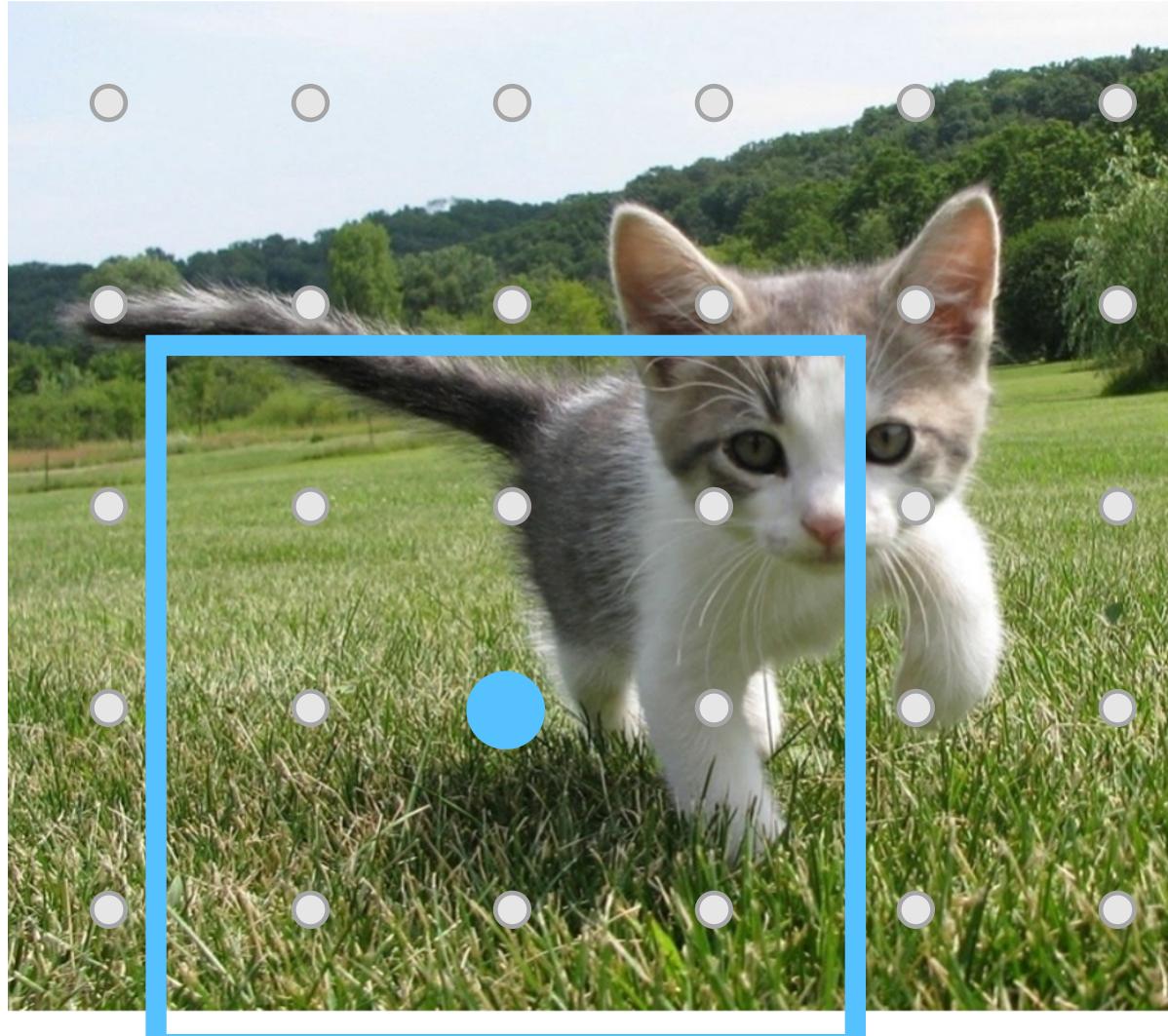


Image features
(e.g. $512 \times 5 \times 6$)

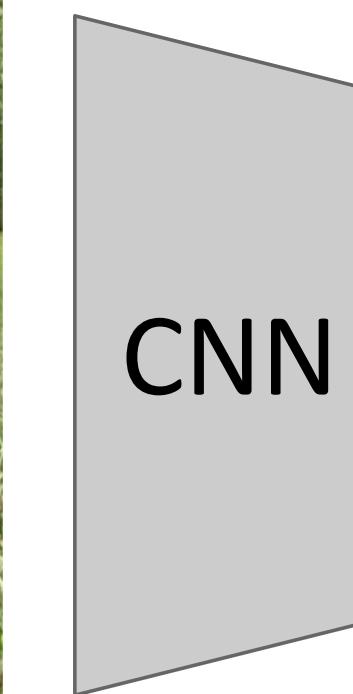
Imagine an **anchor box** of fixed size at each point in the feature map

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. $3 \times 640 \times 480$)



Each feature corresponds to a point in the input

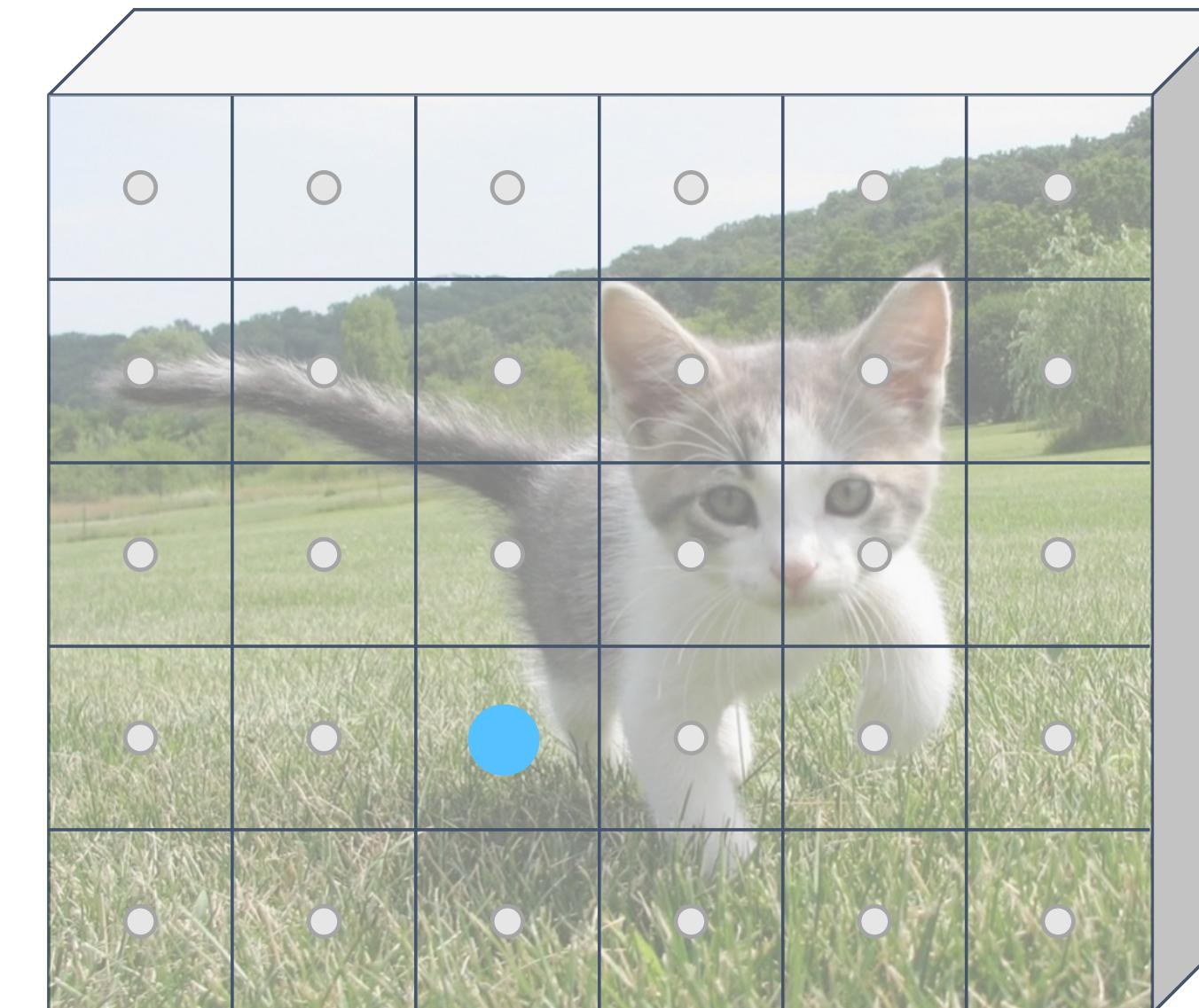
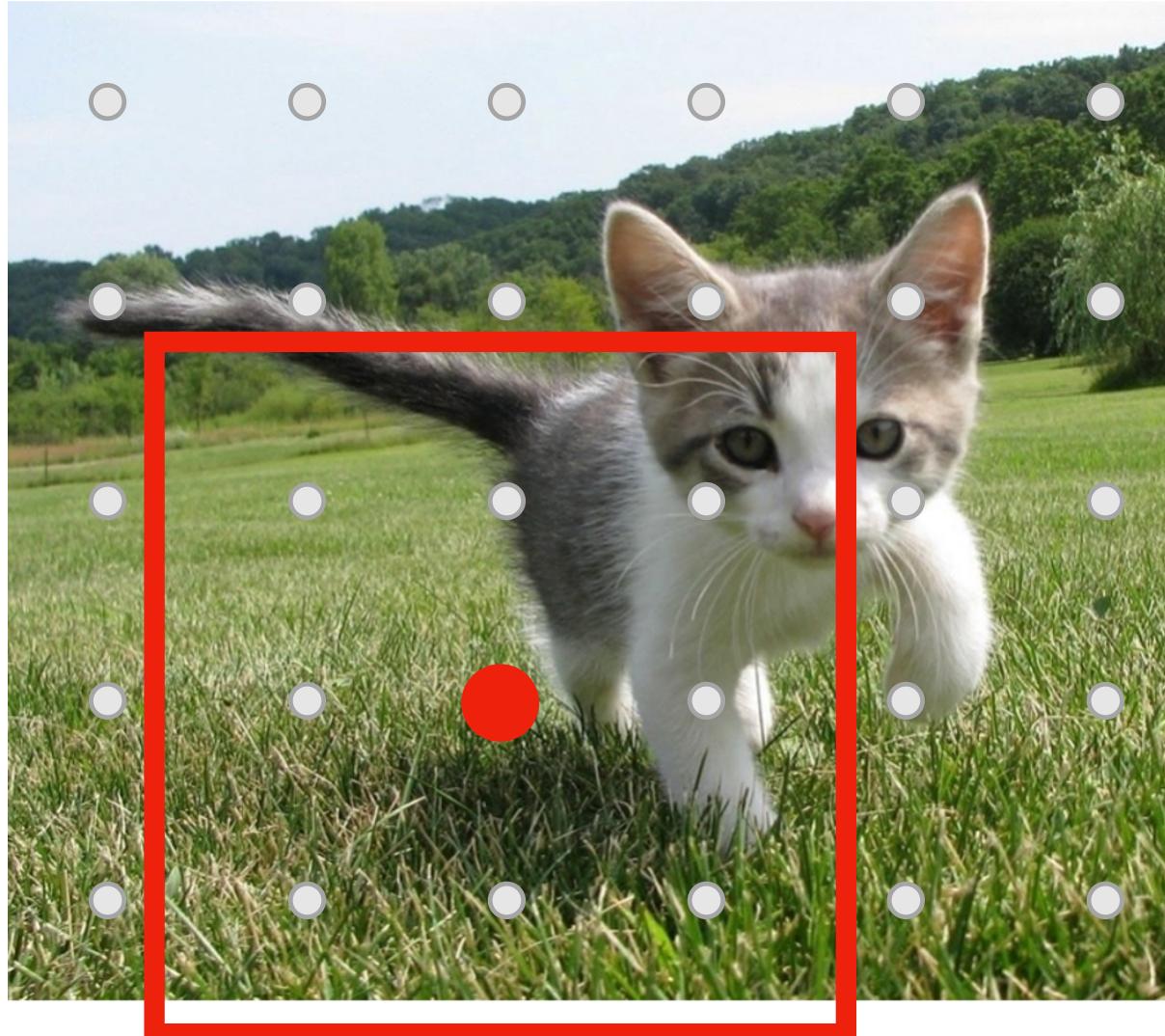


Image features
(e.g. $512 \times 5 \times 6$)

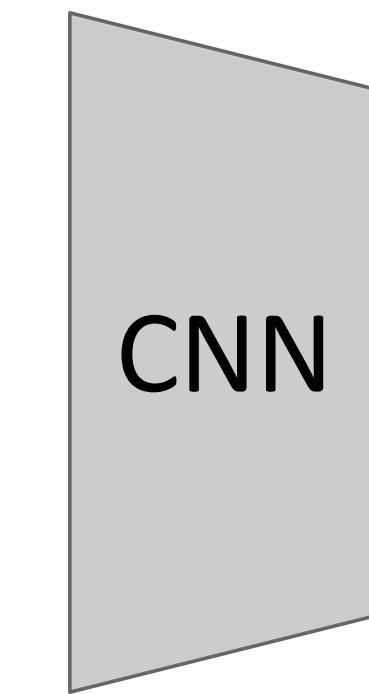
Imagine an **anchor box** of fixed size at each point in the feature map

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. $3 \times 640 \times 480$)



Each feature corresponds to a point in the input

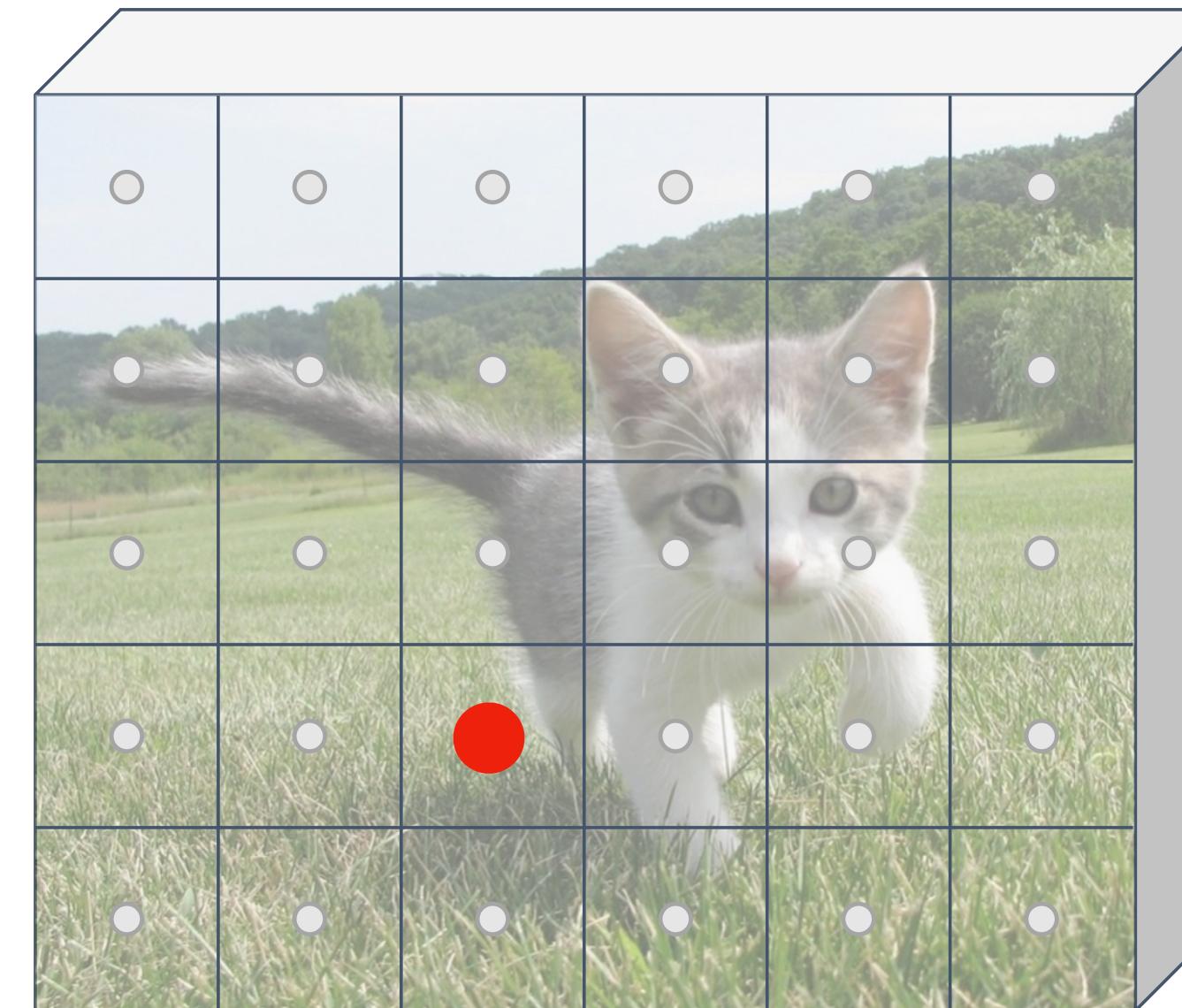


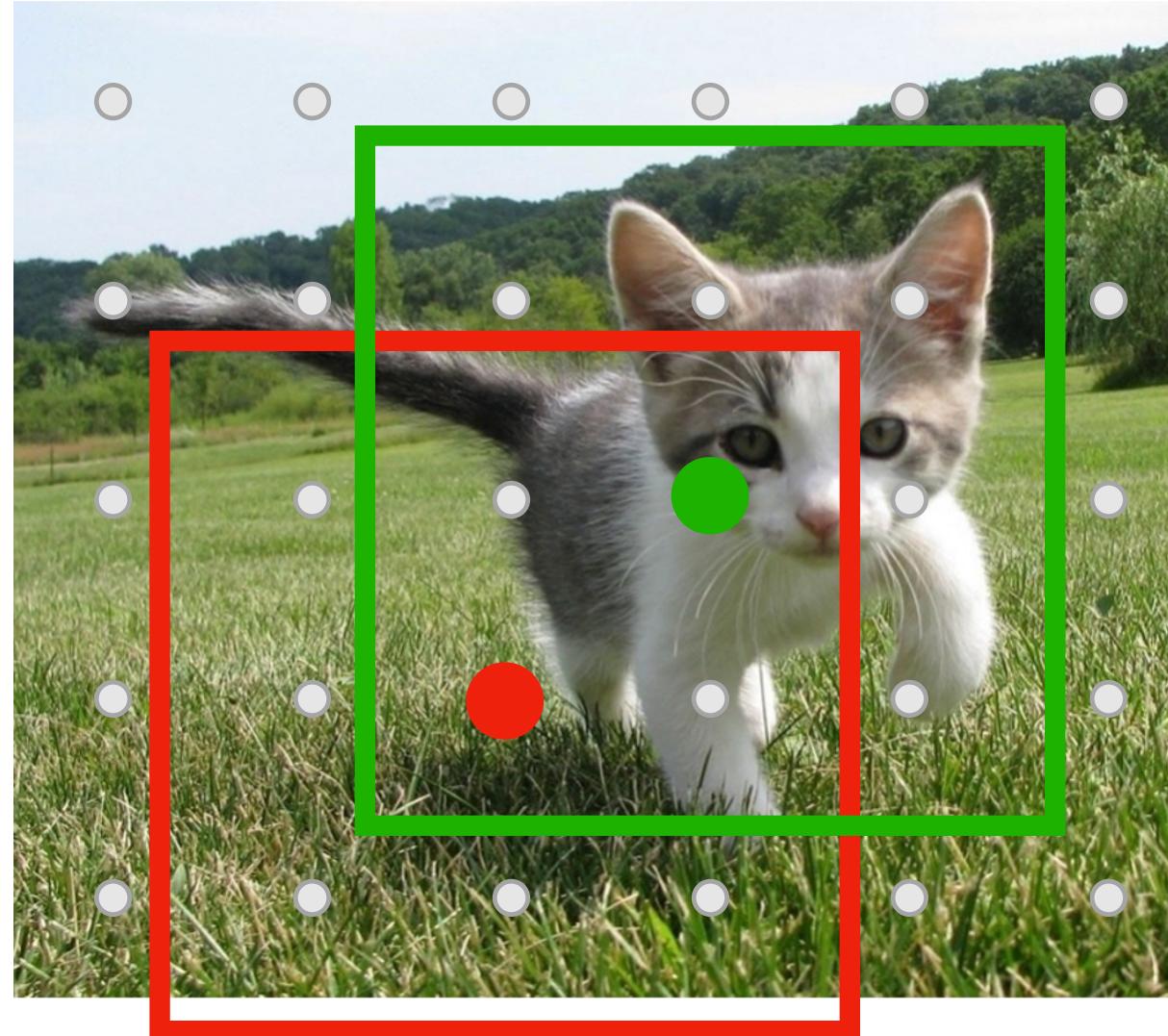
Image features
(e.g. $512 \times 5 \times 6$)

Imagine an **anchor box** of fixed size at each point in the feature map

Classify each anchor as positive (object) or negative (no object)

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. $3 \times 640 \times 480$)

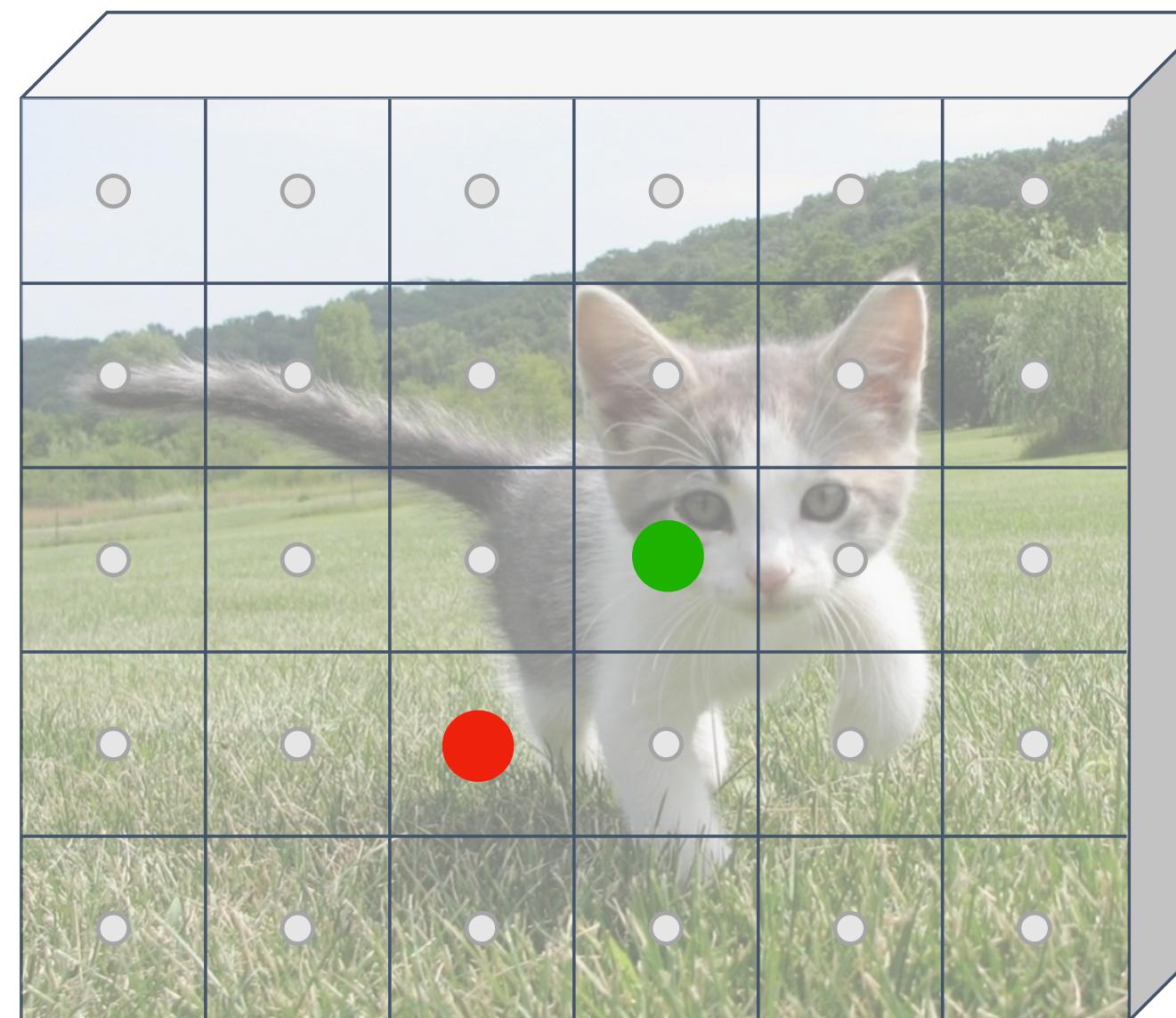
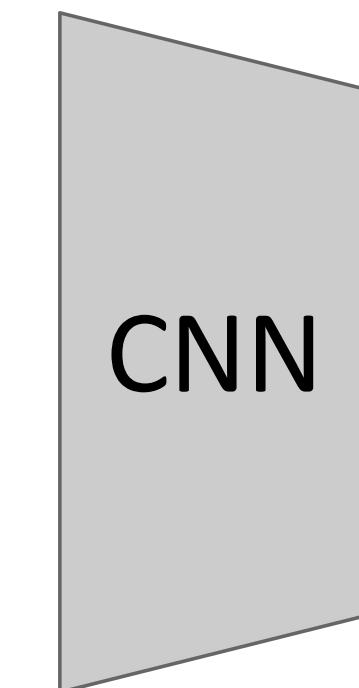


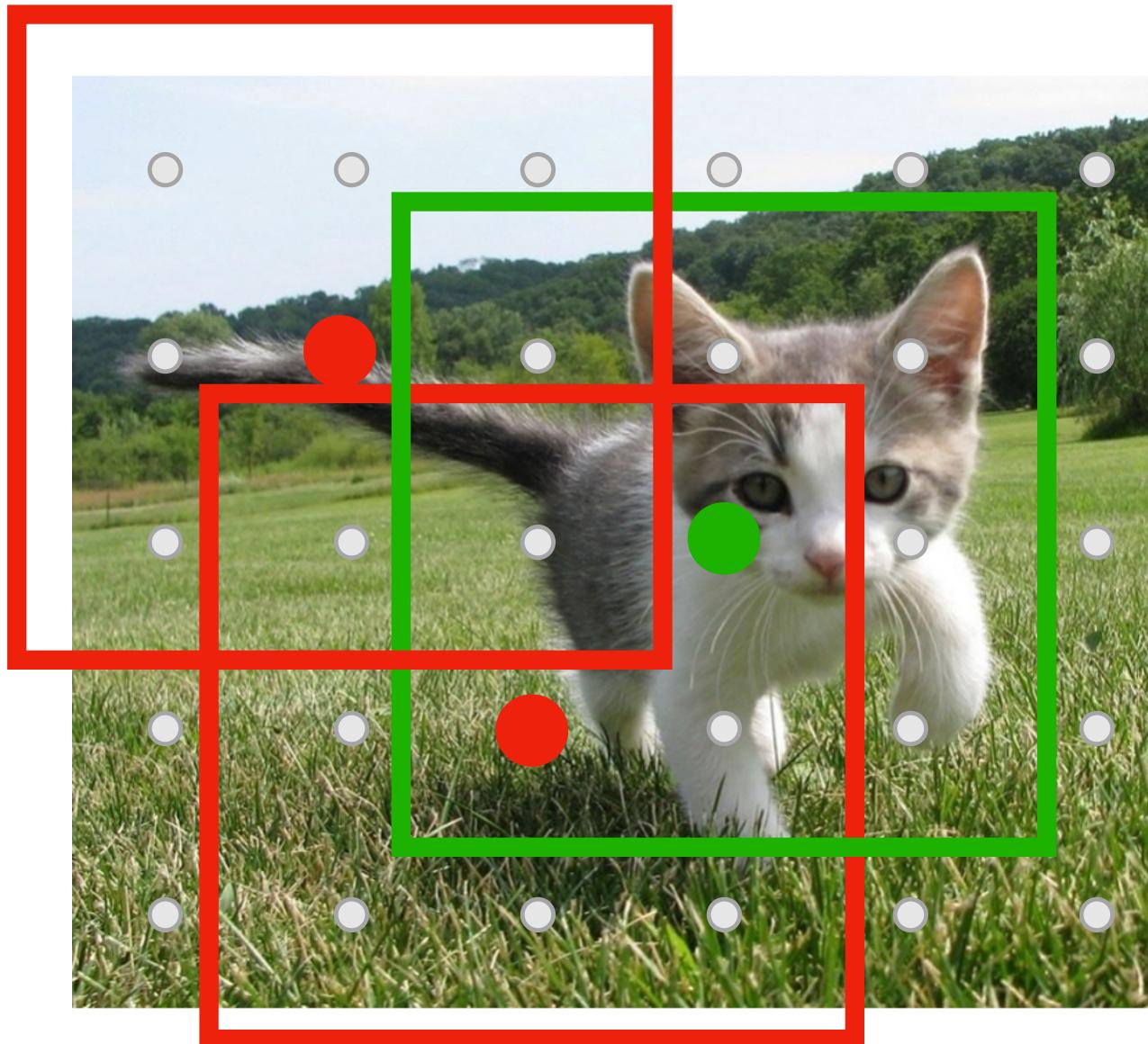
Image features
(e.g. $512 \times 5 \times 6$)

Imagine an **anchor box** of fixed size at each point in the feature map

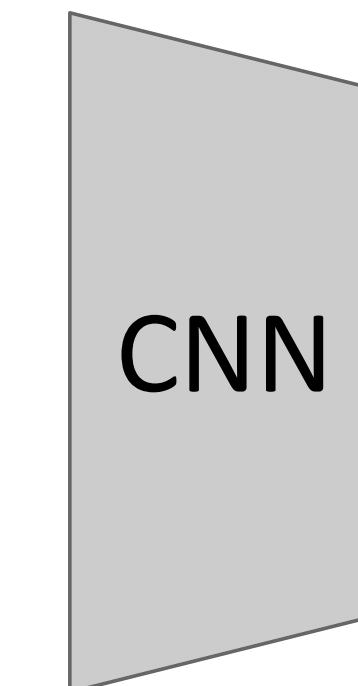
Classify each anchor as positive (object) or negative (no object)

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. $3 \times 640 \times 480$)



Each feature corresponds to a point in the input

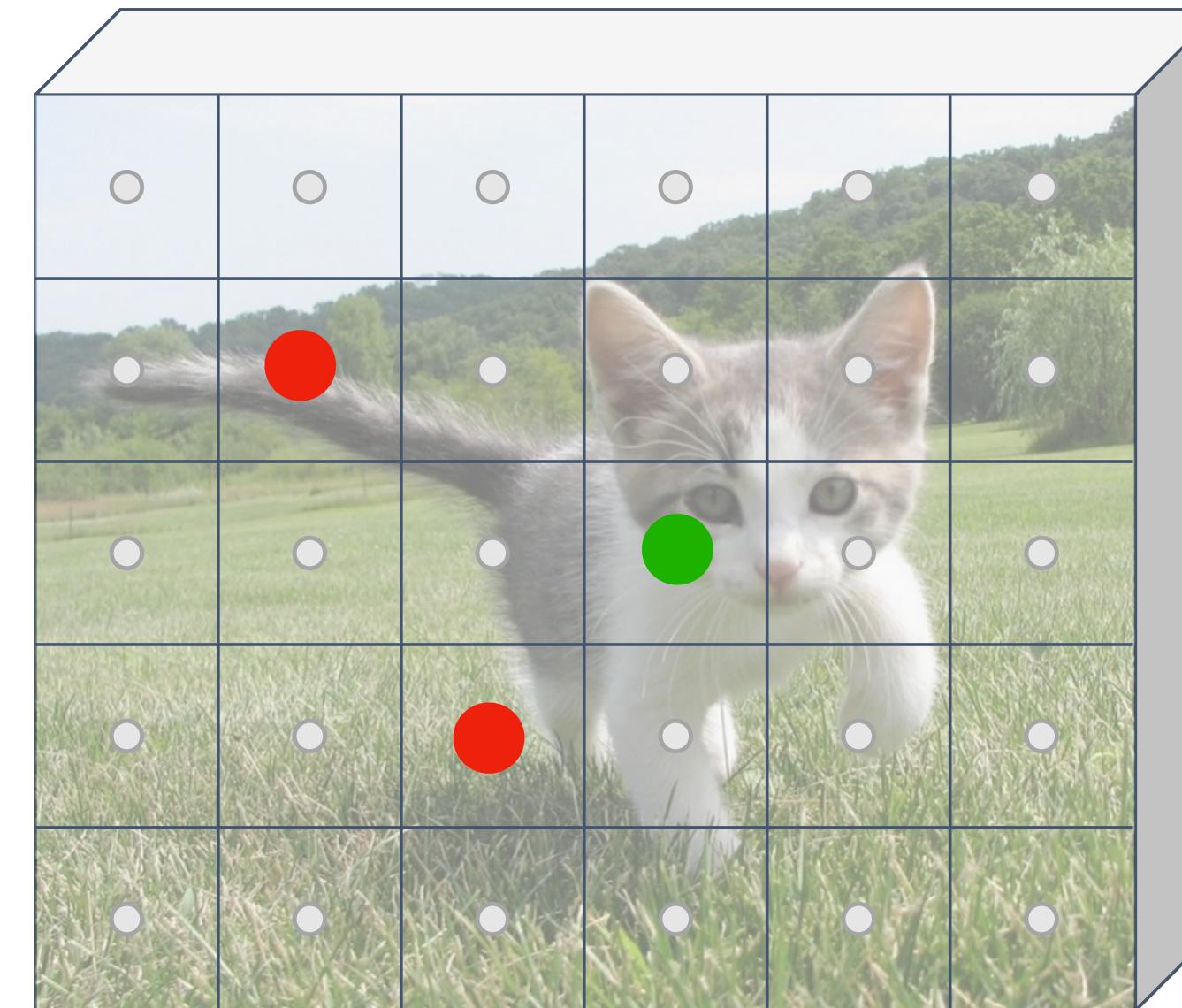


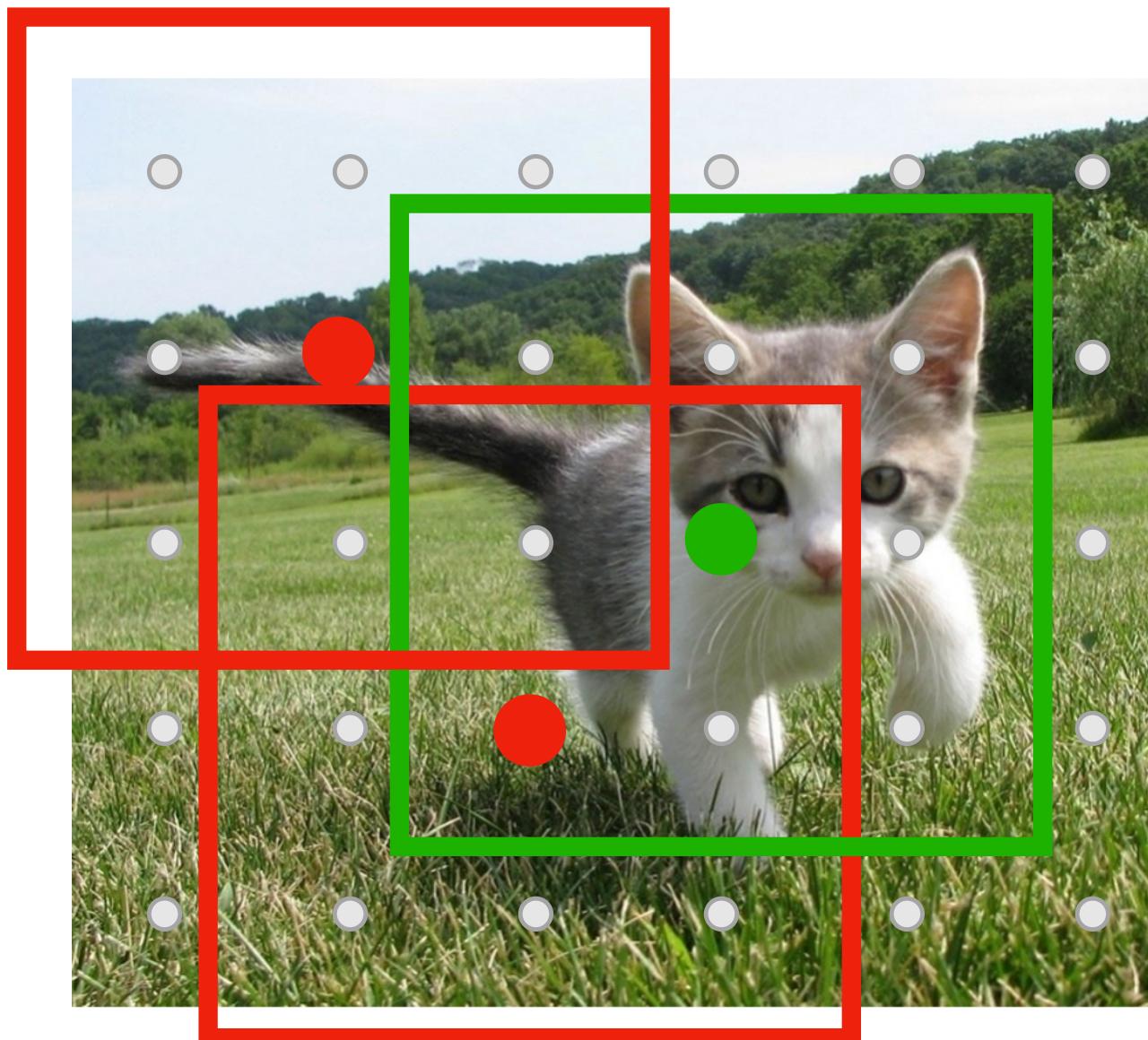
Image features
(e.g. $512 \times 5 \times 6$)

Imagine an **anchor box** of fixed size at each point in the feature map

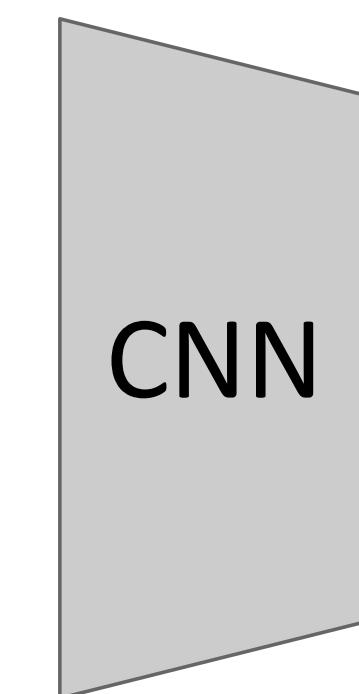
Classify each anchor as positive (object) or negative (no object)

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. $3 \times 640 \times 480$)



Each feature corresponds to a point in the input

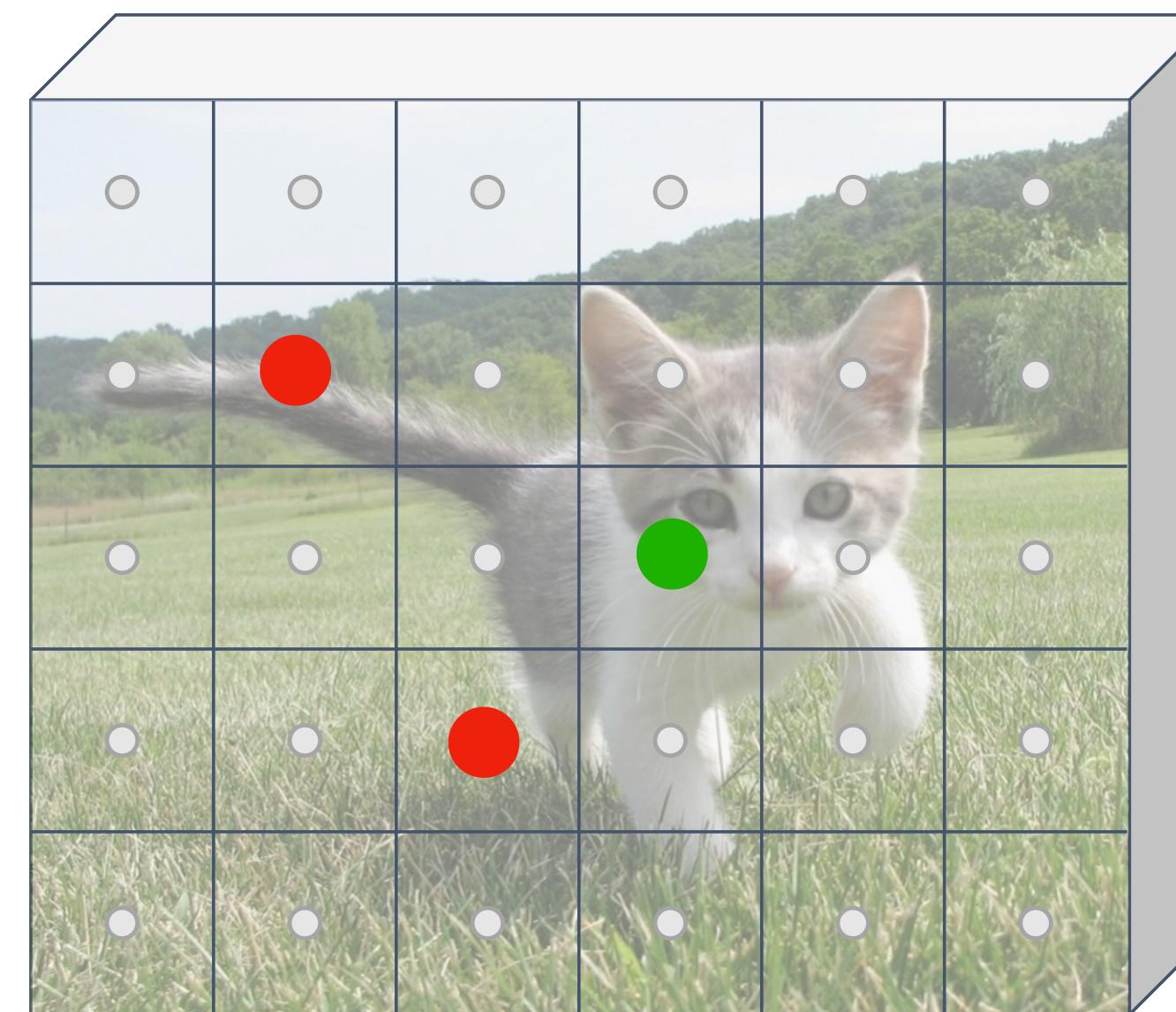


Image features
(e.g. $512 \times 5 \times 6$)

Predict object vs not object scores for all anchors with a conv layer (512 input filters, 2 output filters)

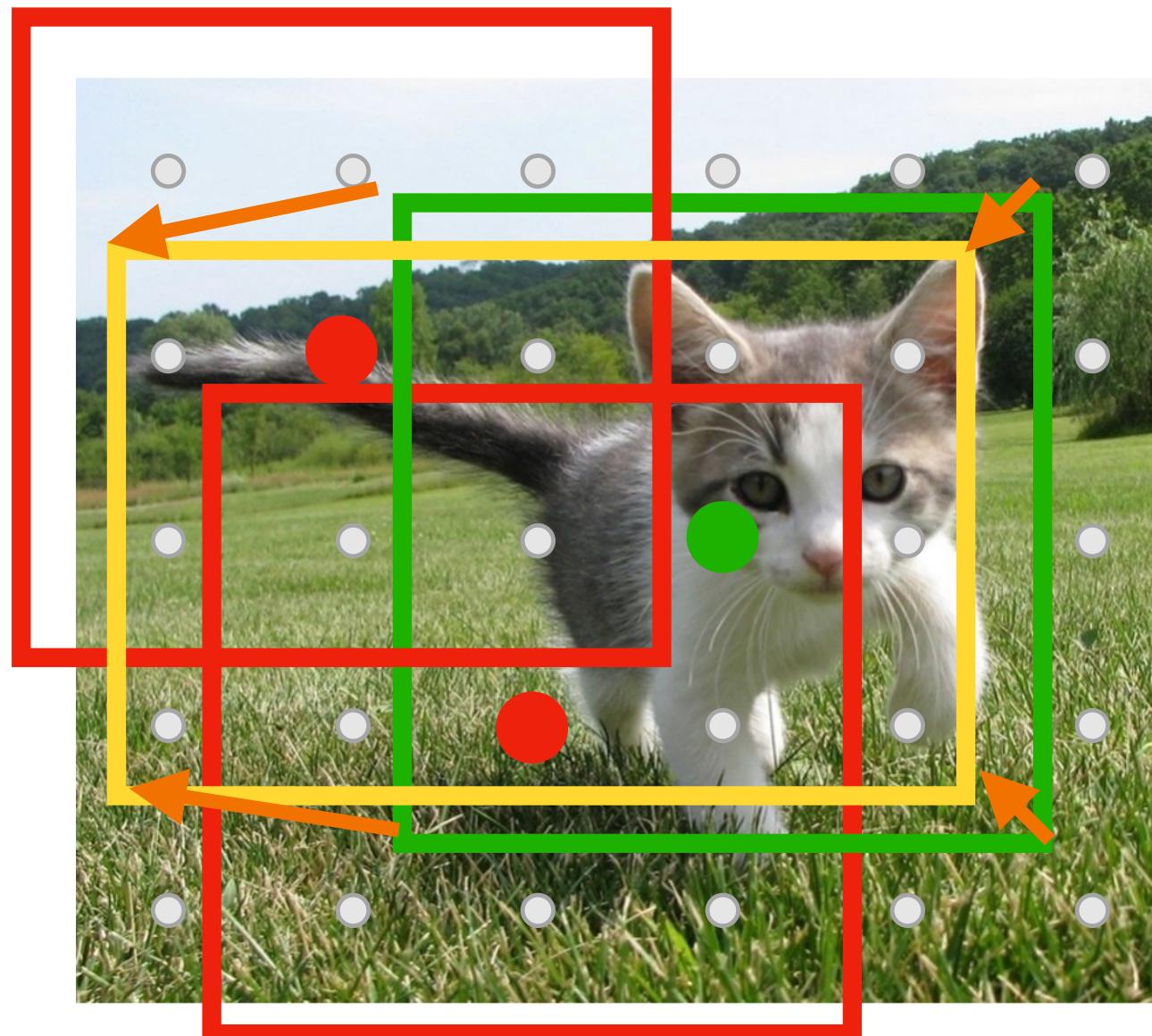


Anchor is object?
 $2 \times 5 \times 6$

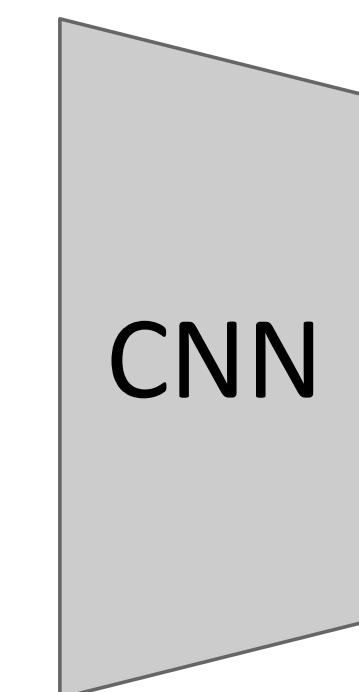
Classify each anchor as positive (object) or negative (no object)

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. $3 \times 640 \times 480$)



Each feature corresponds to a point in the input

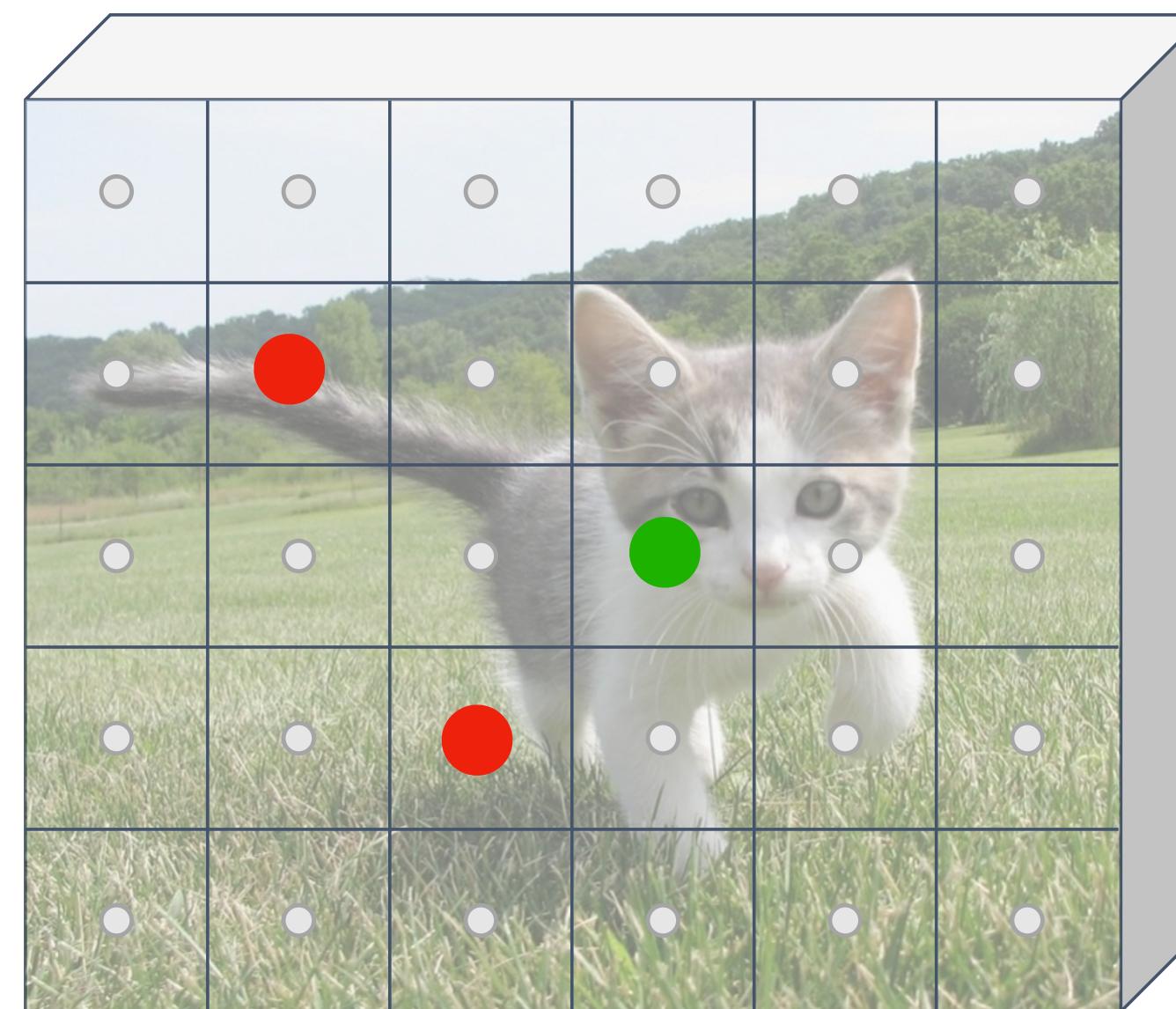


Image features
(e.g. $512 \times 5 \times 6$)

For positive anchors, also predict a transform that converting the anchor to the GT box (like R-CNN)

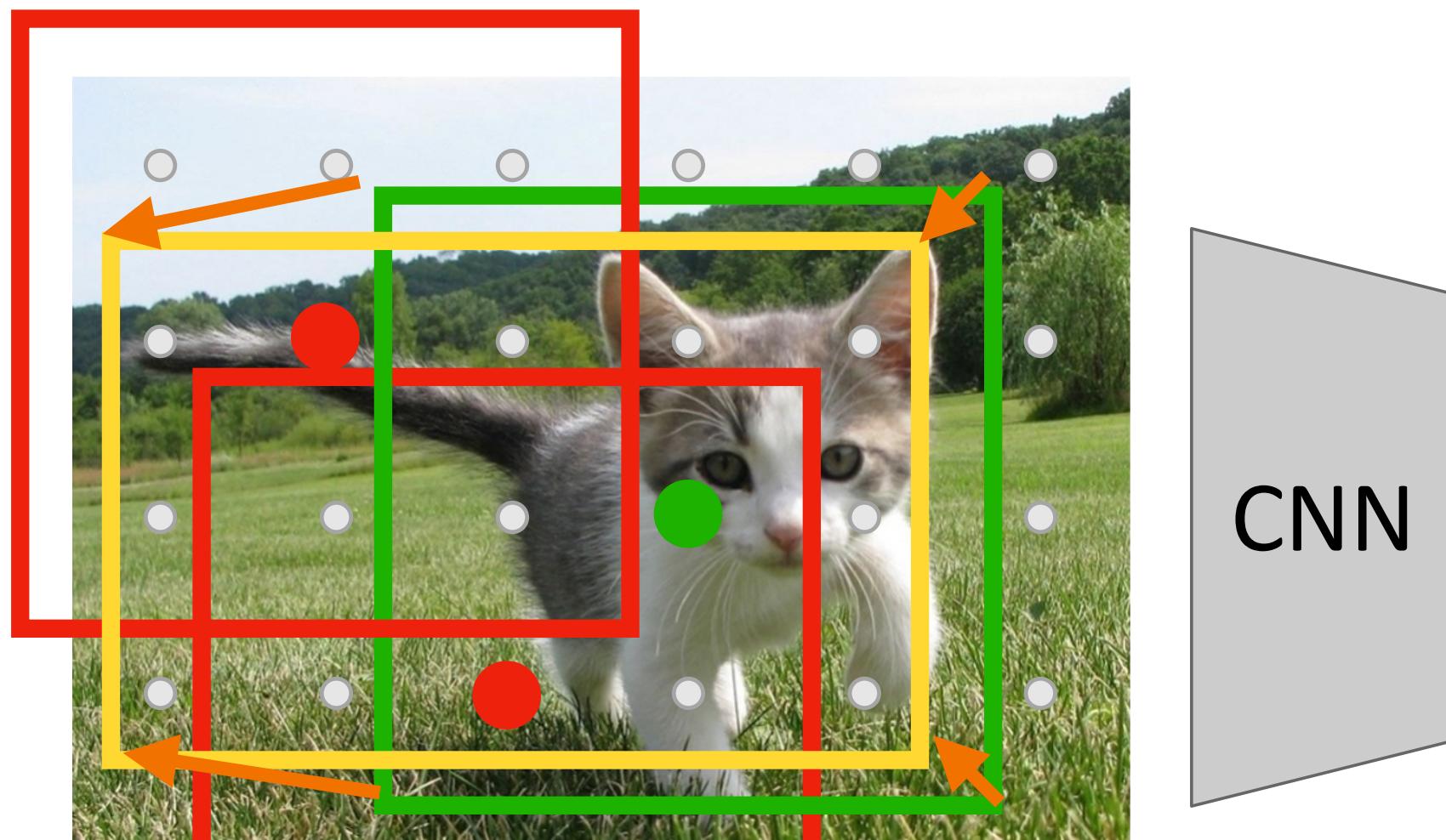
Anchor is object?
 $2 \times 5 \times 6$



Classify each anchor as positive (object) or negative (no object)

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. $3 \times 640 \times 480$)

Each feature corresponds to a point in the input

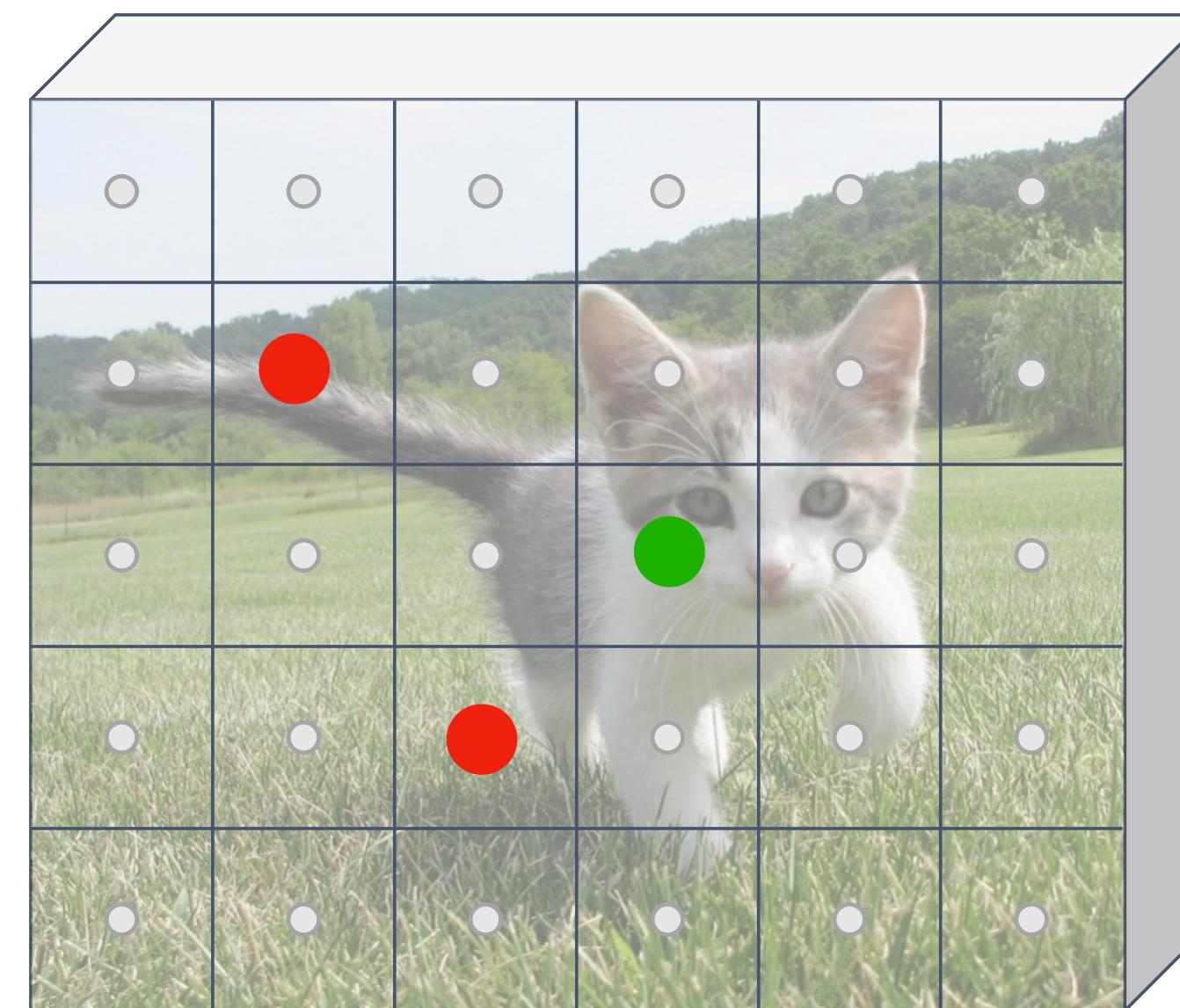
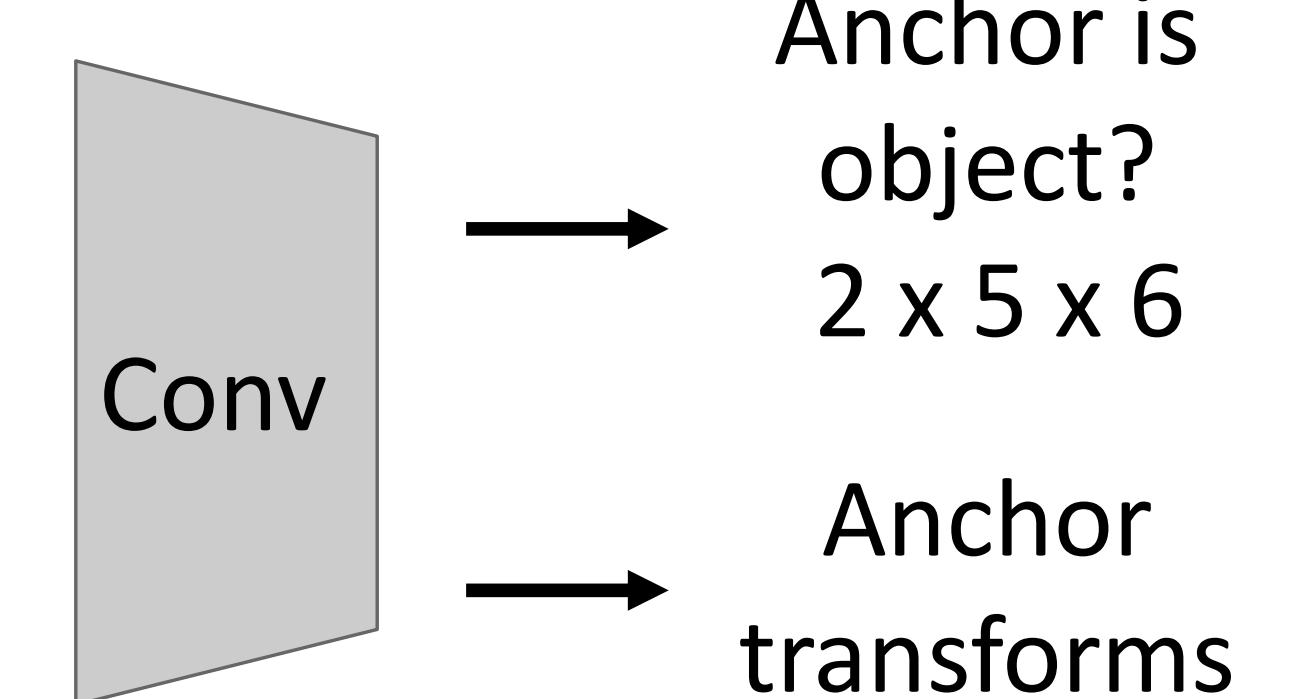


Image features
(e.g. $512 \times 5 \times 6$)

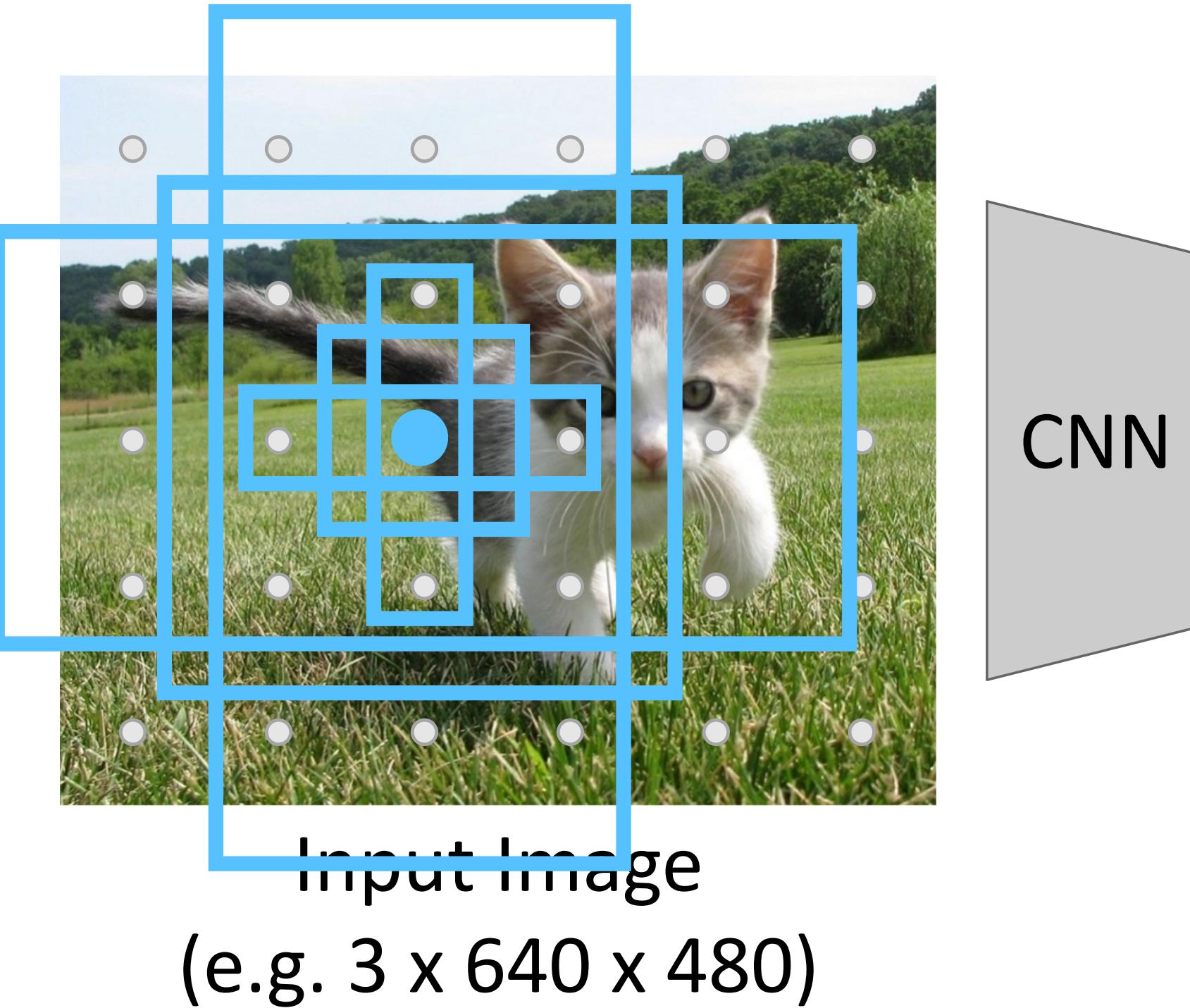
For **positive anchors**, also predict a **transform** that converting the anchor to the **GT box** (like R-CNN)



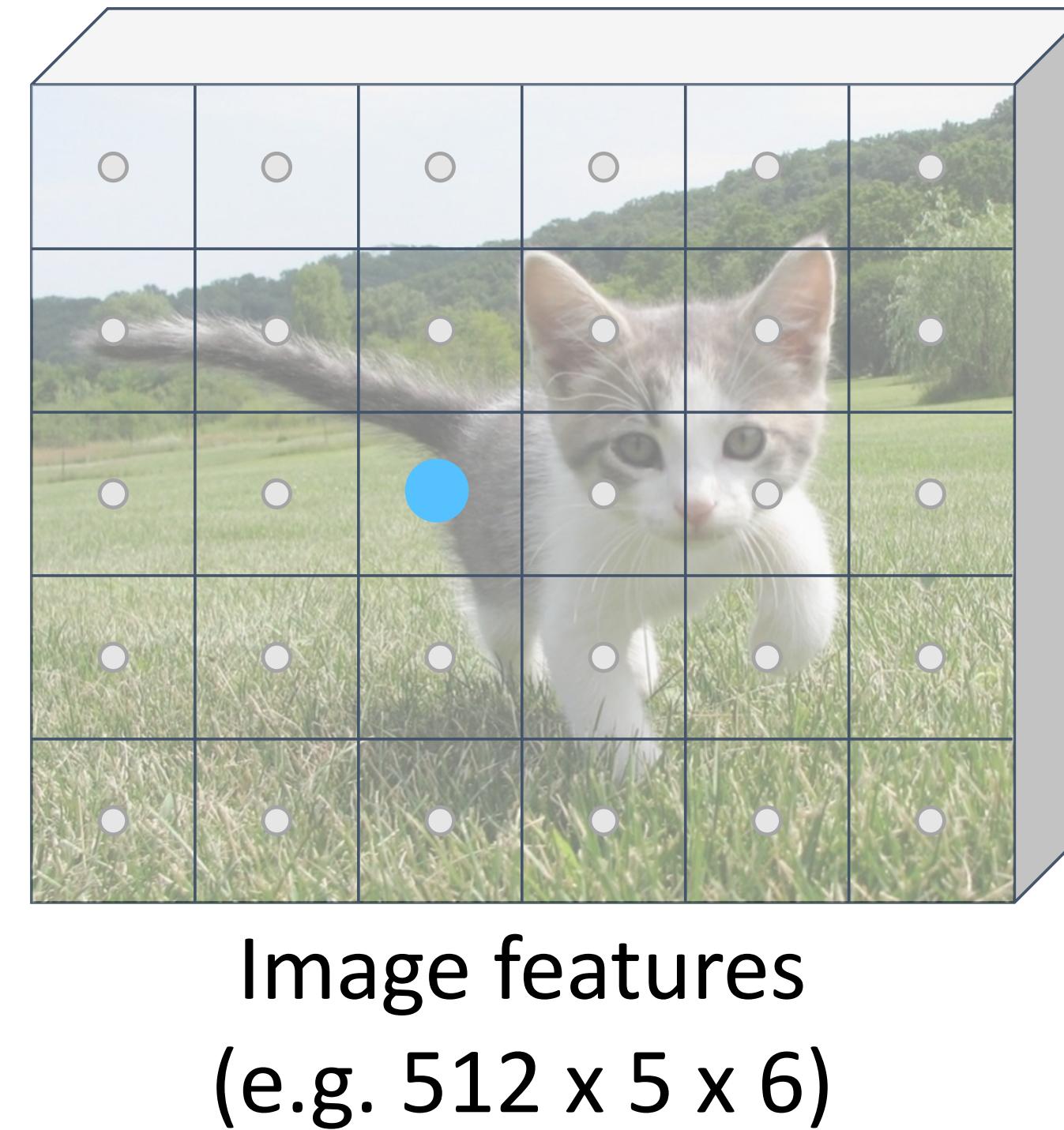
Classify each anchor as **positive (object)** or **negative (no object)**

Region Proposal Network (RPN)

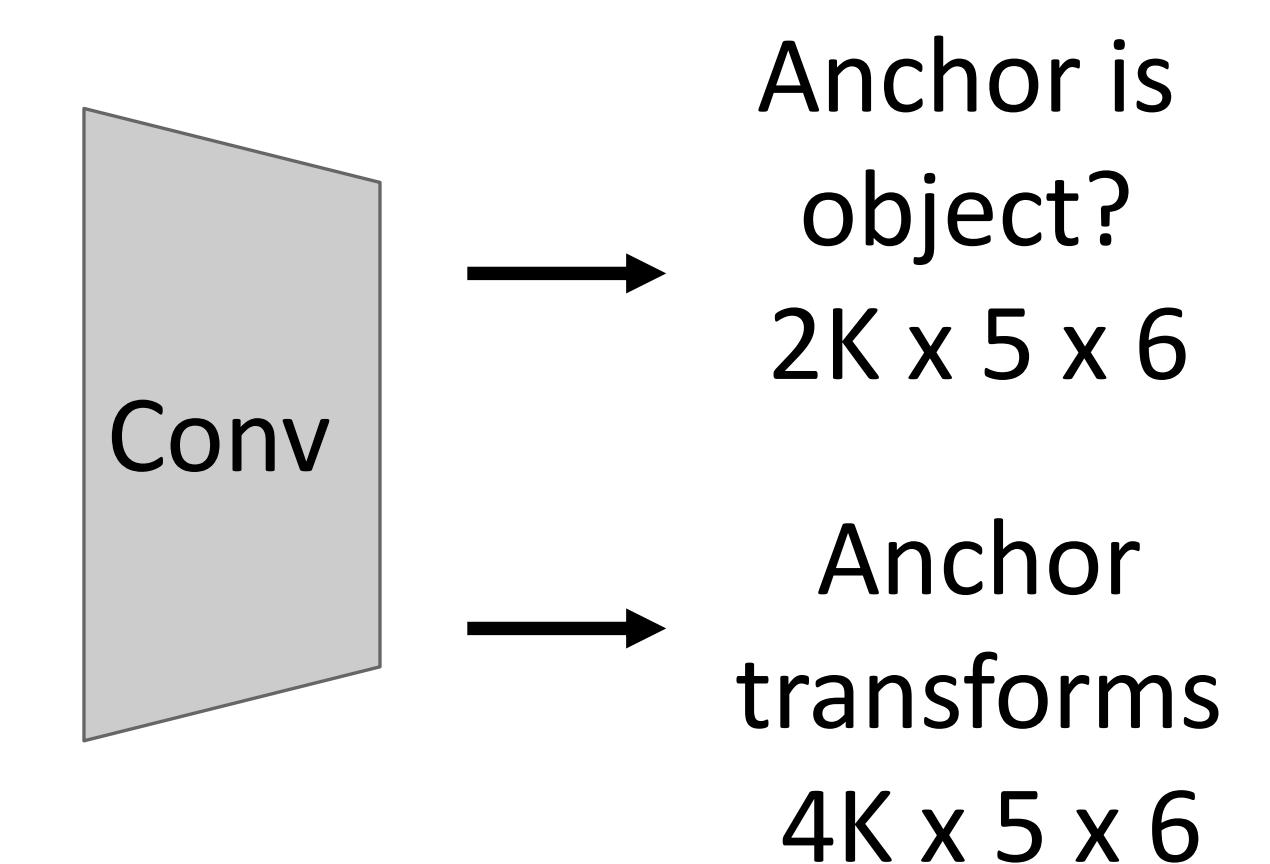
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input

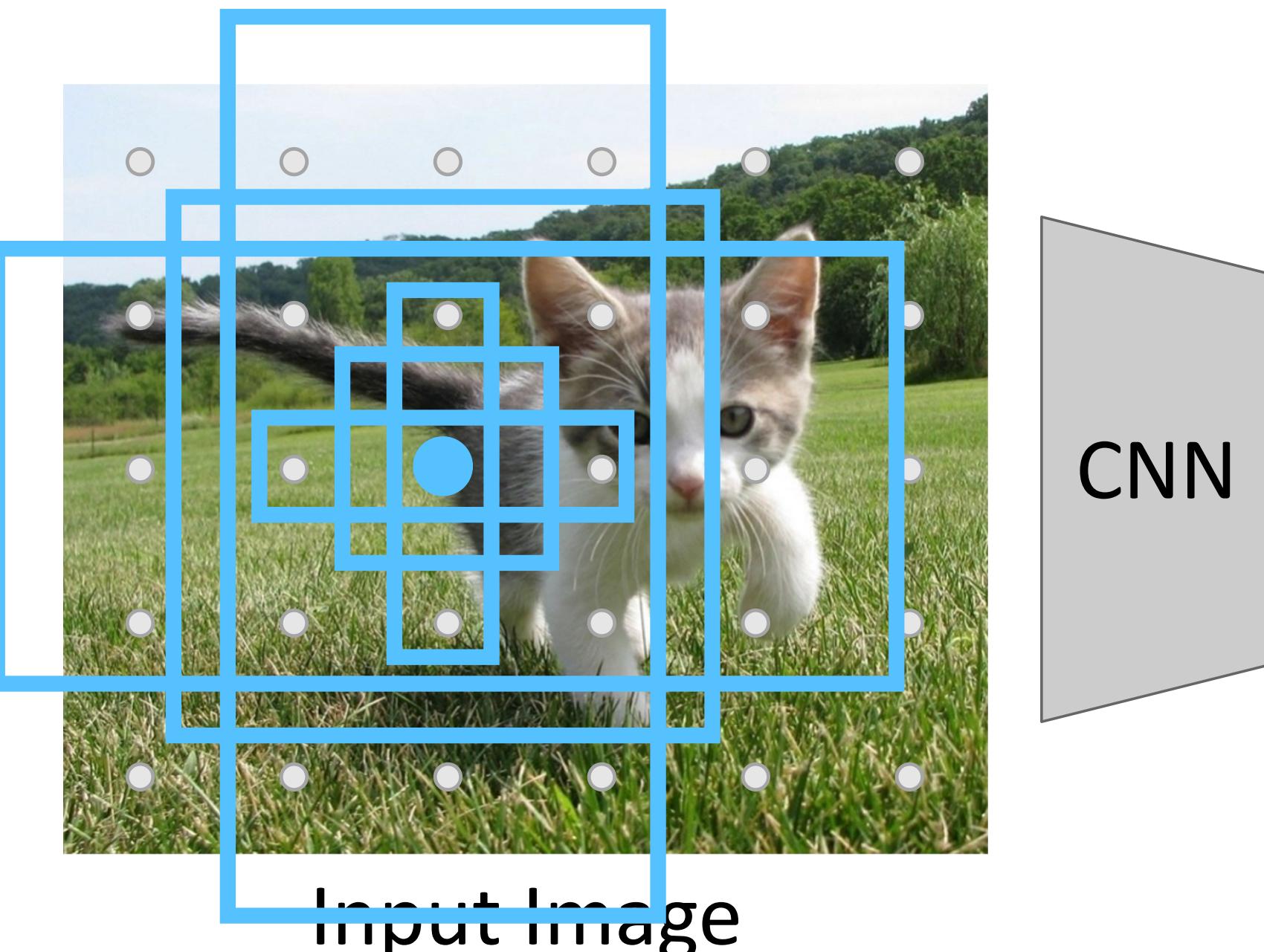


In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here $K = 6$)

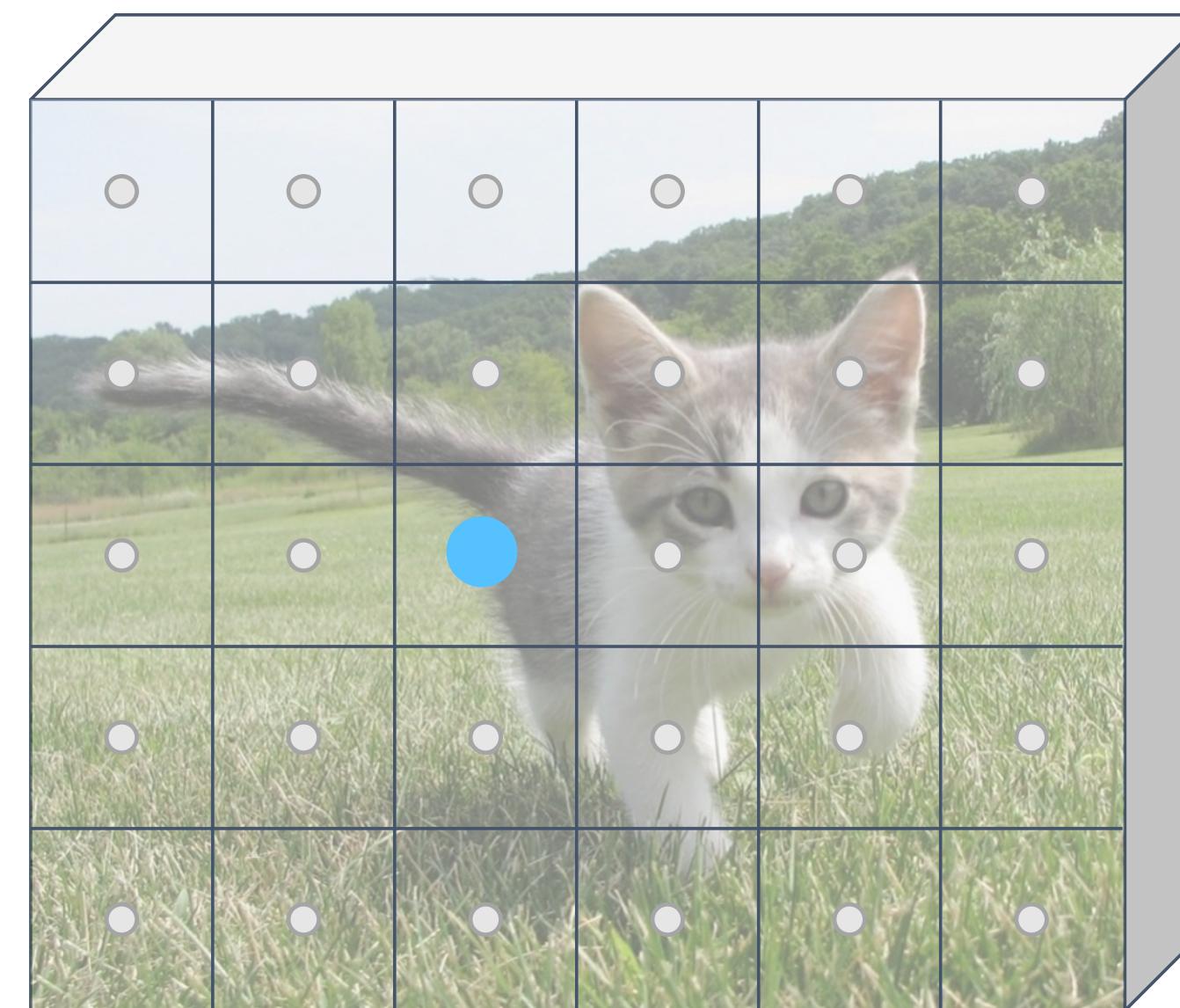


Region Proposal Network (RPN)

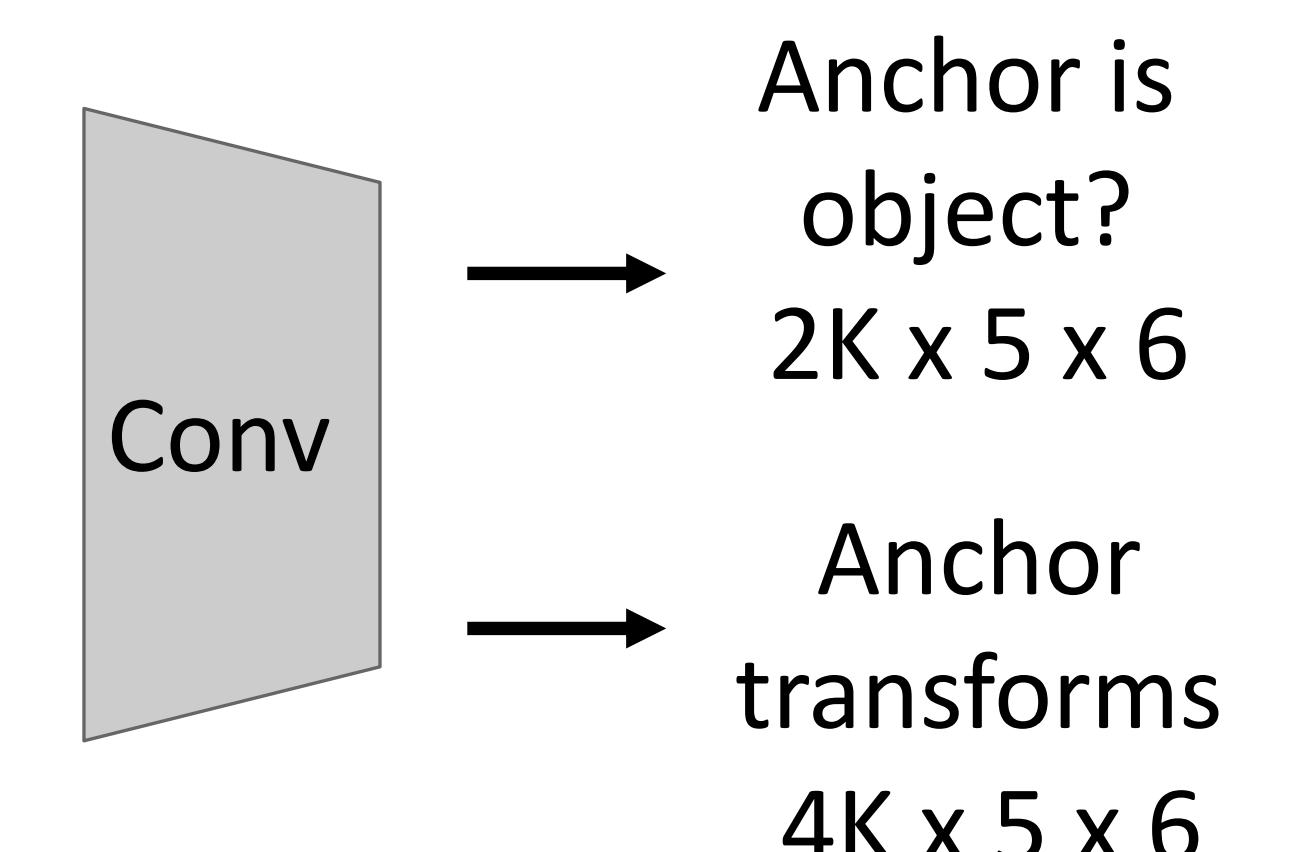
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



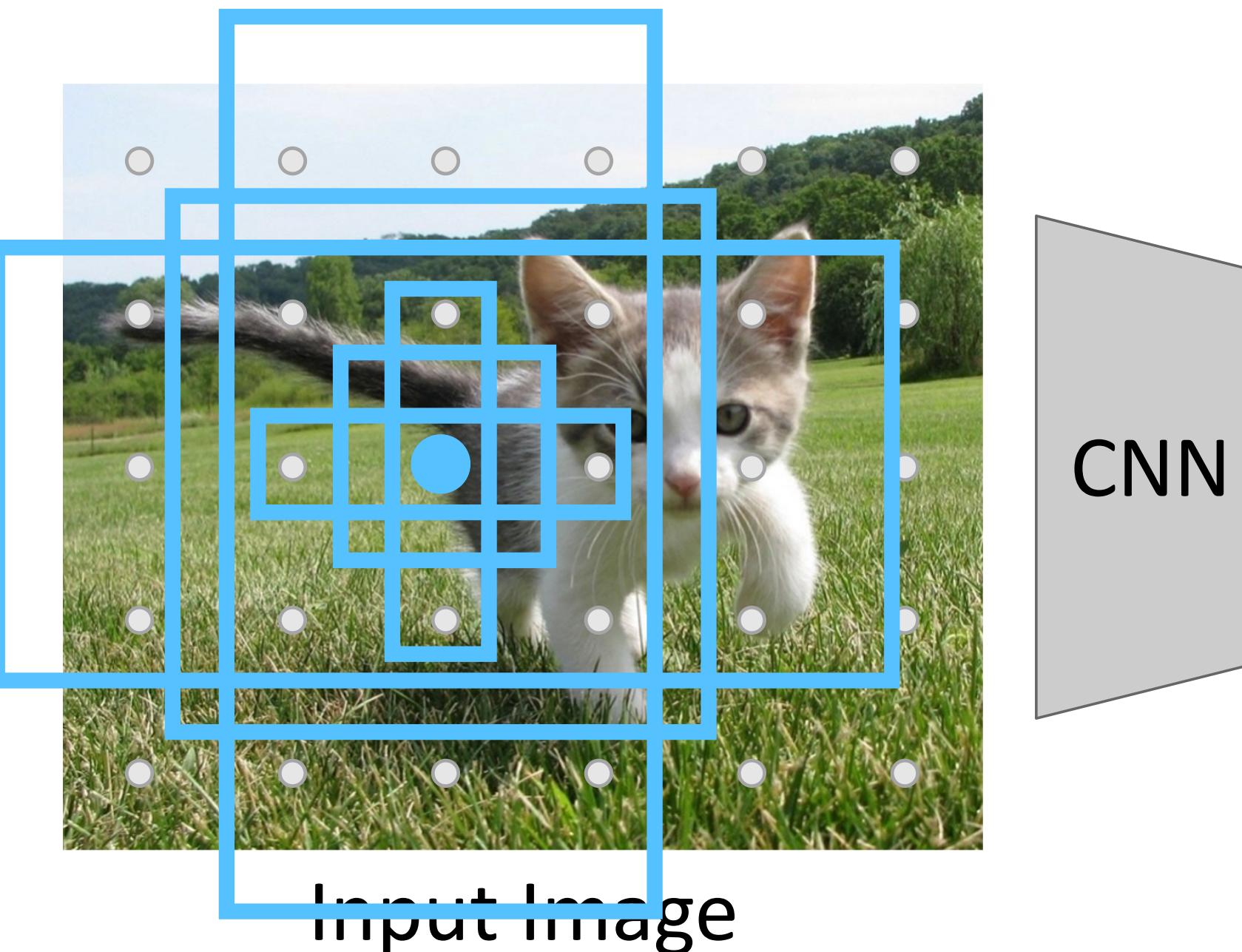
In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here $K = 6$)



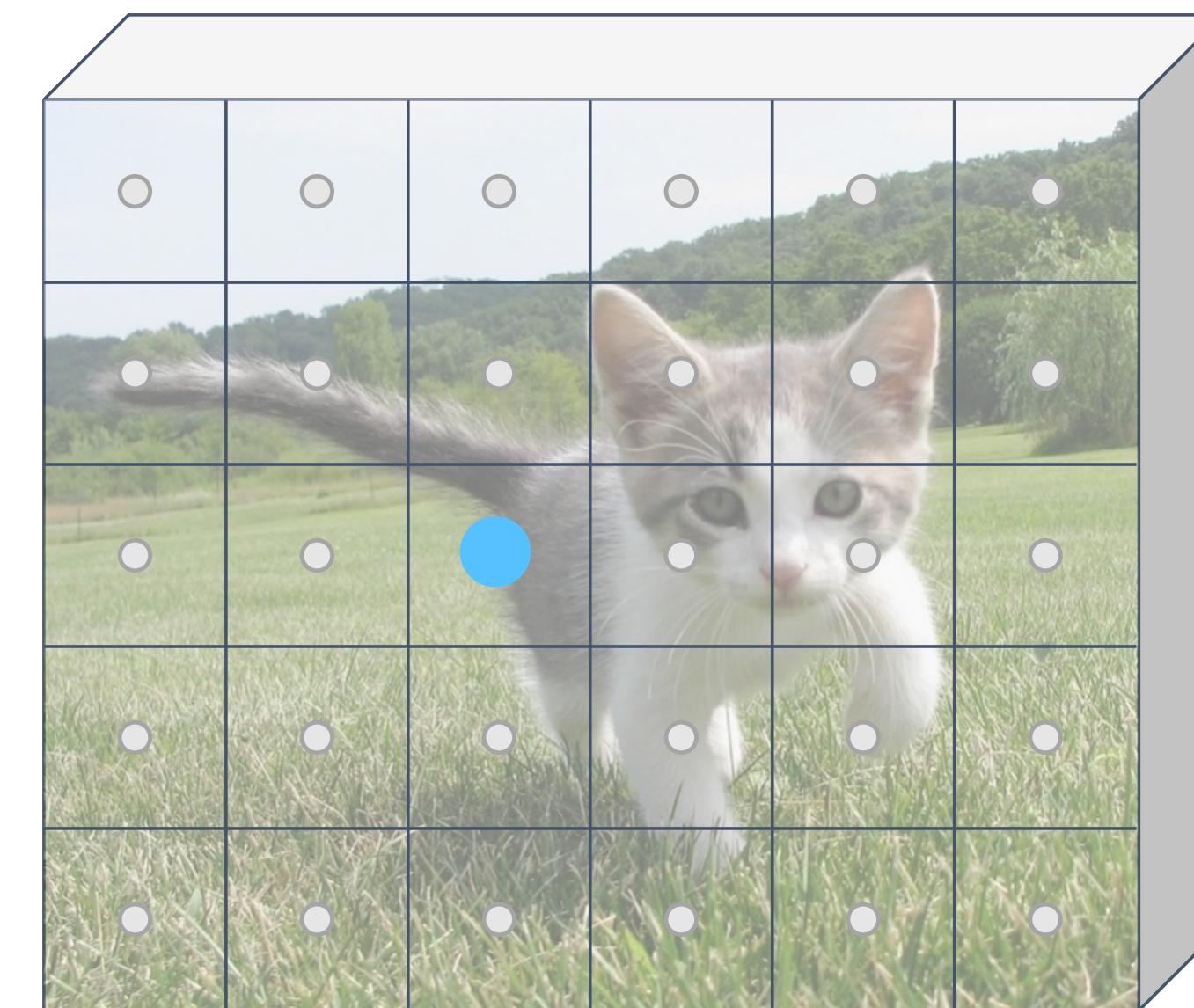
During training, supervised positive / negative anchors and box transforms like R-CNN

Region Proposal Network (RPN)

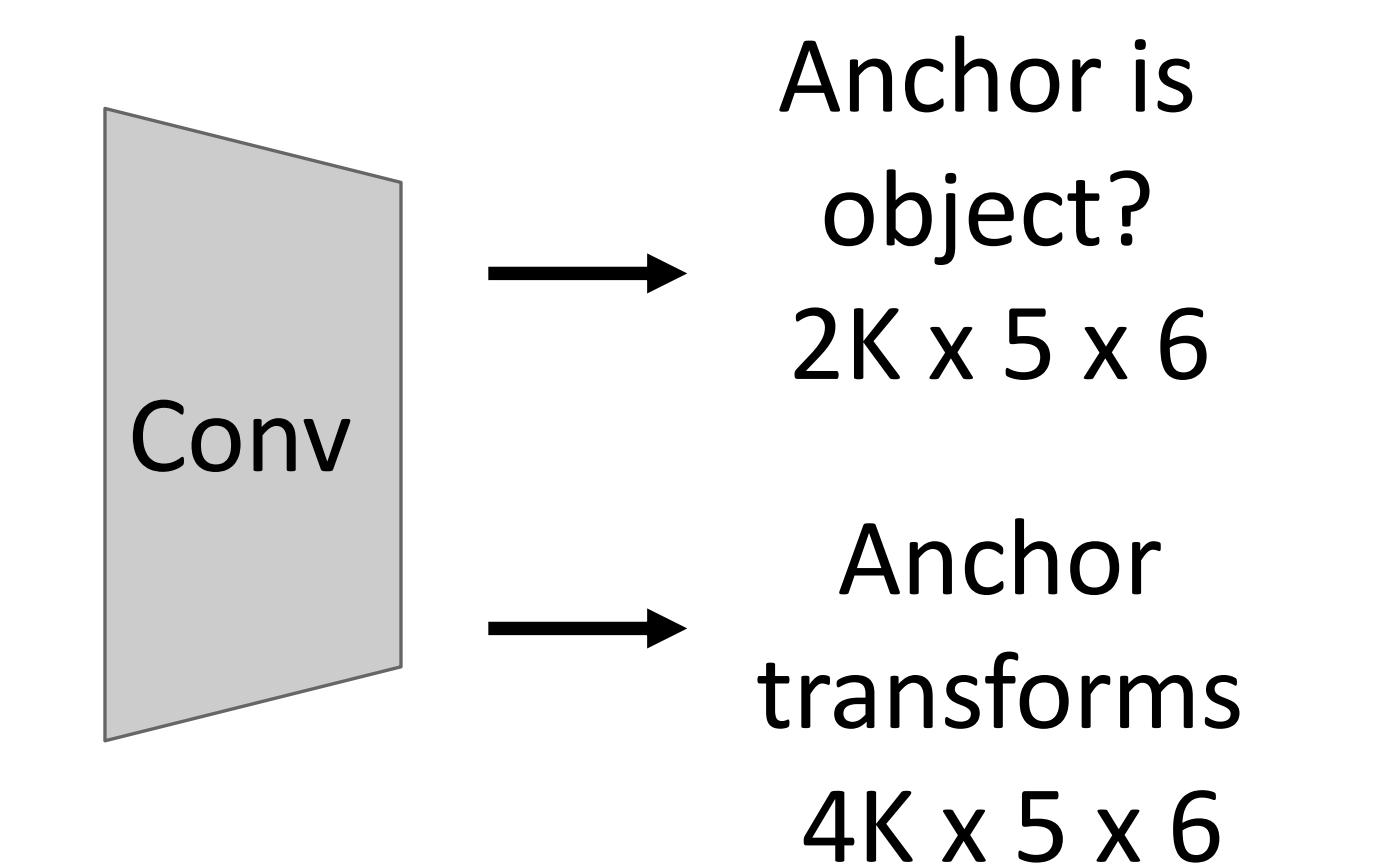
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



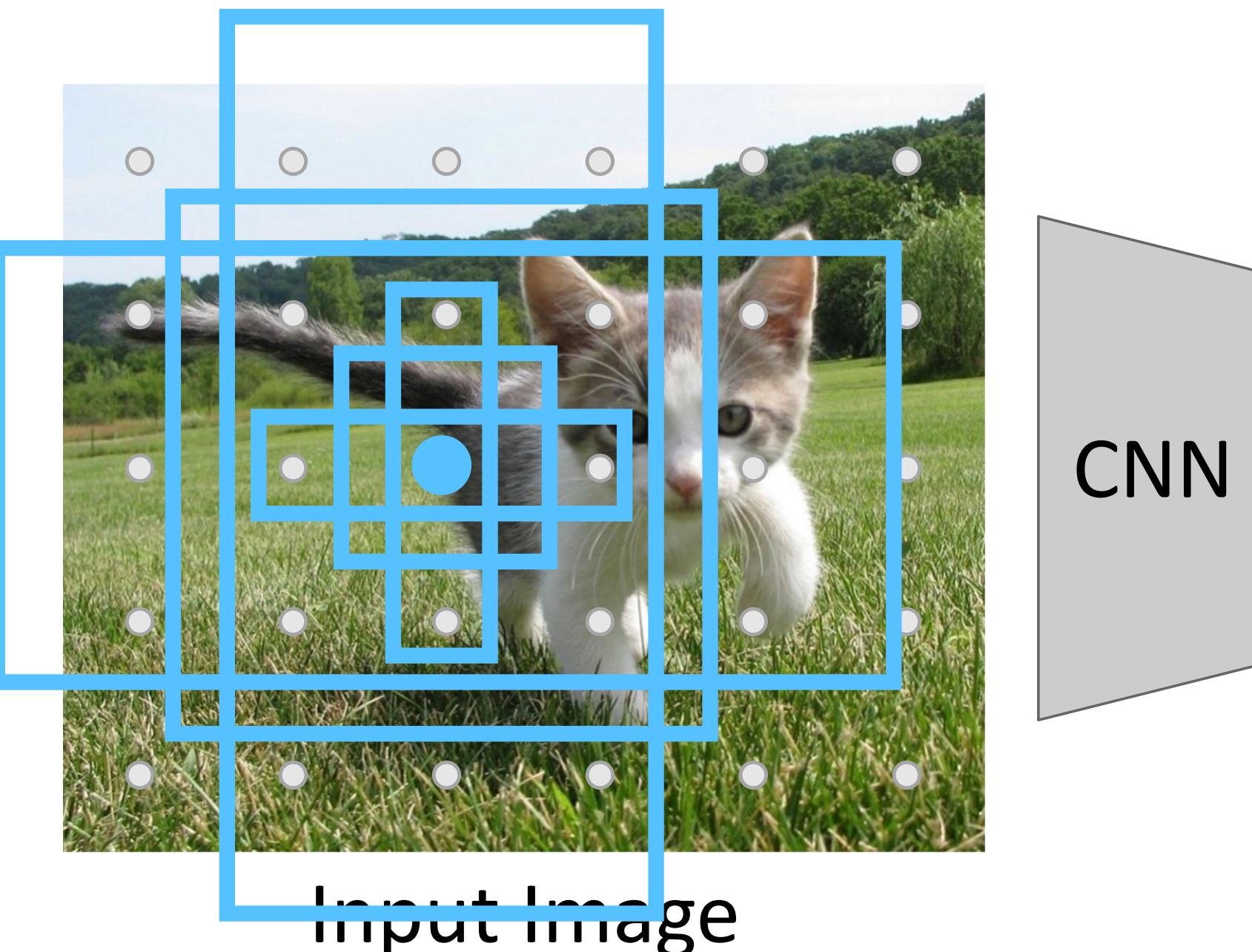
In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here $K = 6$)



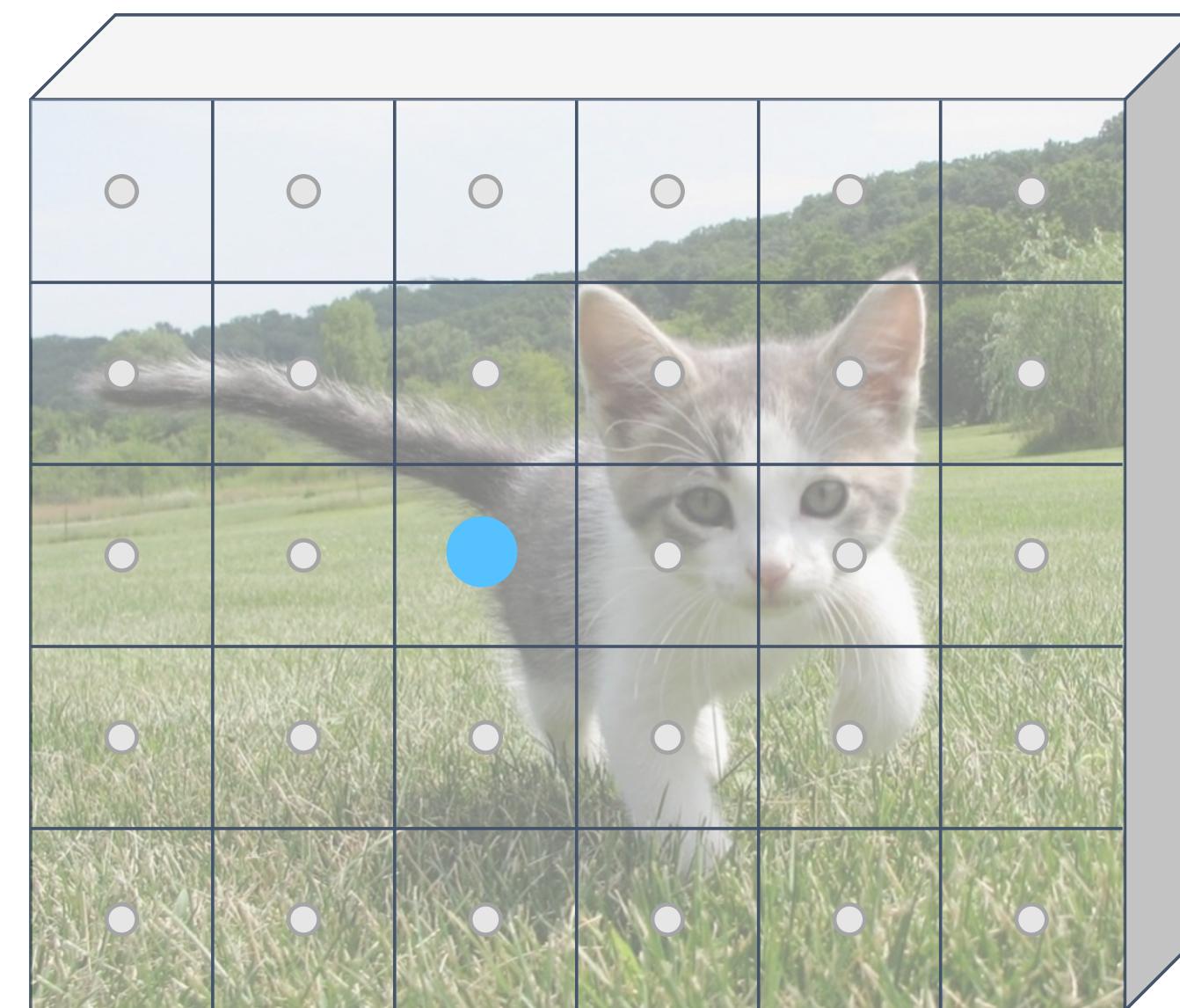
Positive anchors: ≥ 0.7 IoU with some GT box (plus highest IoU to each GT)

Region Proposal Network (RPN)

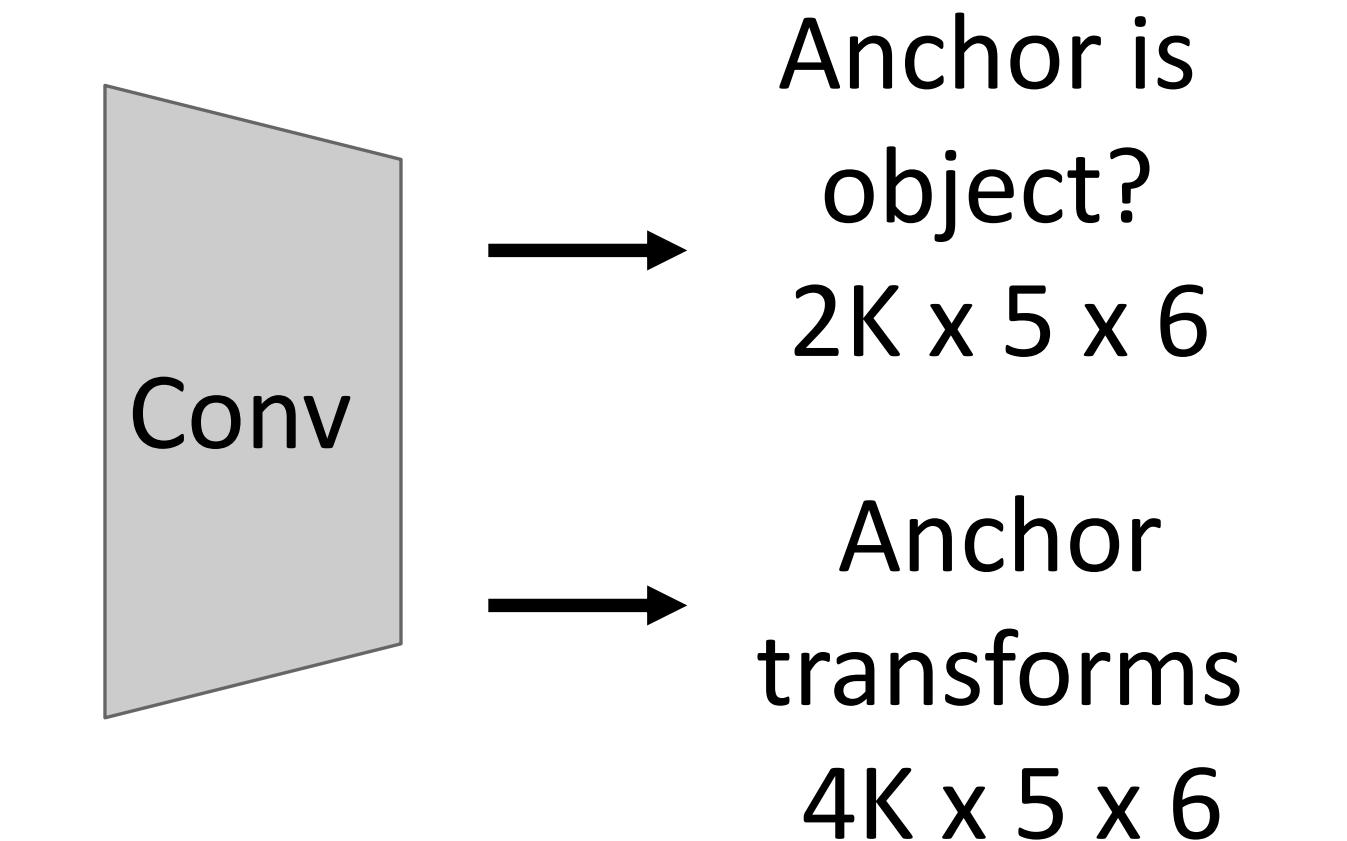
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



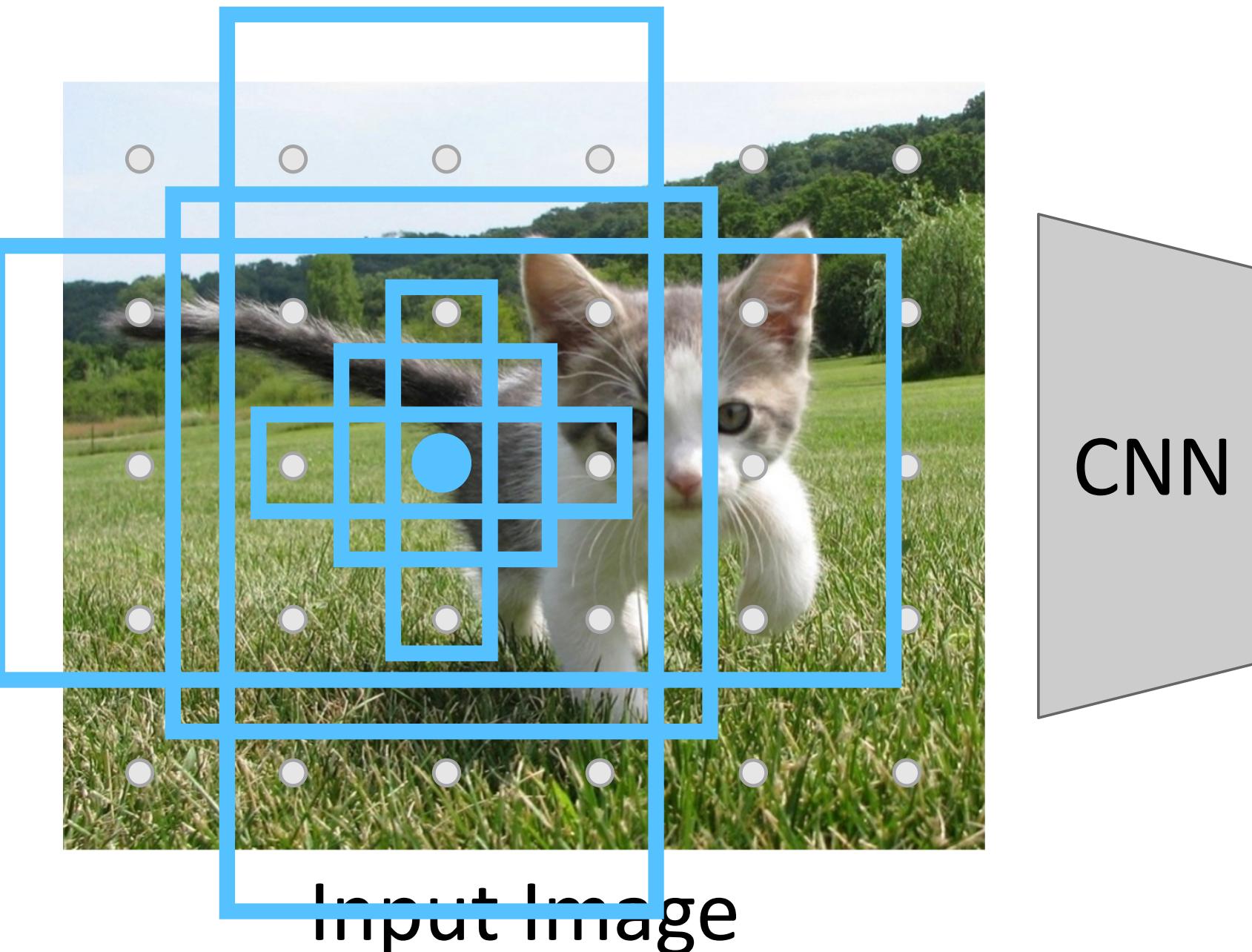
In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here $K = 6$)



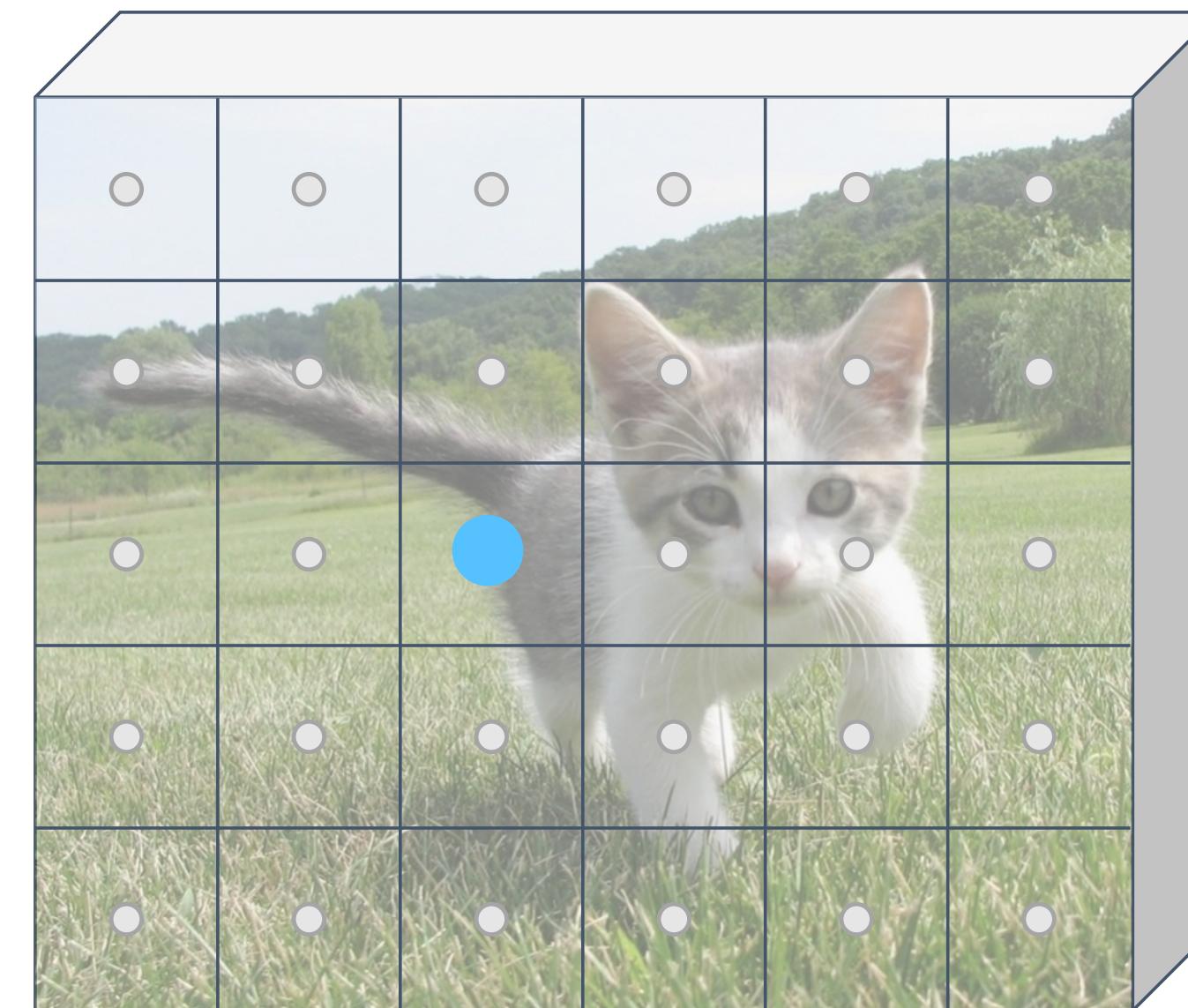
Negative anchors: < 0.3 IoU with all GT boxes. Don't supervise transforms for negative boxes.

Region Proposal Network (RPN)

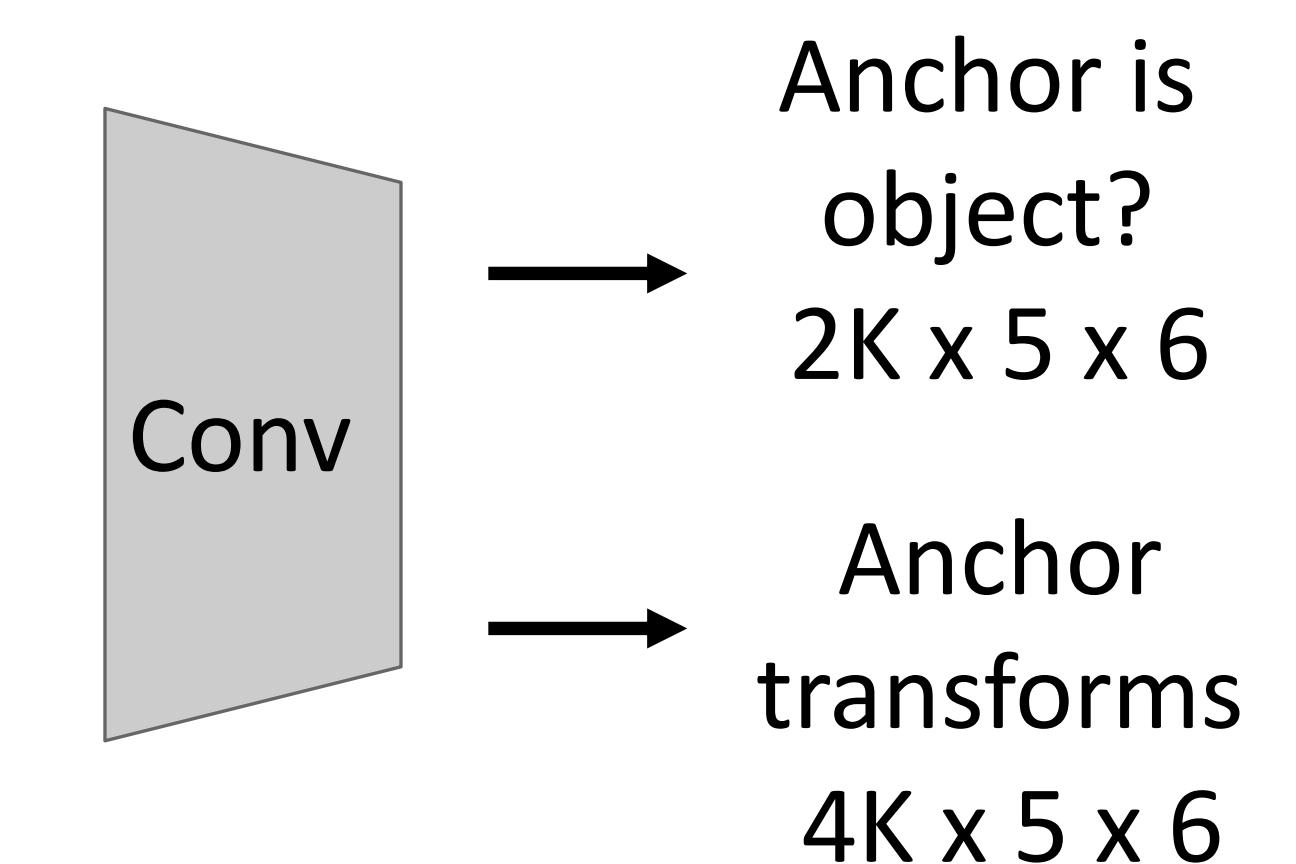
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



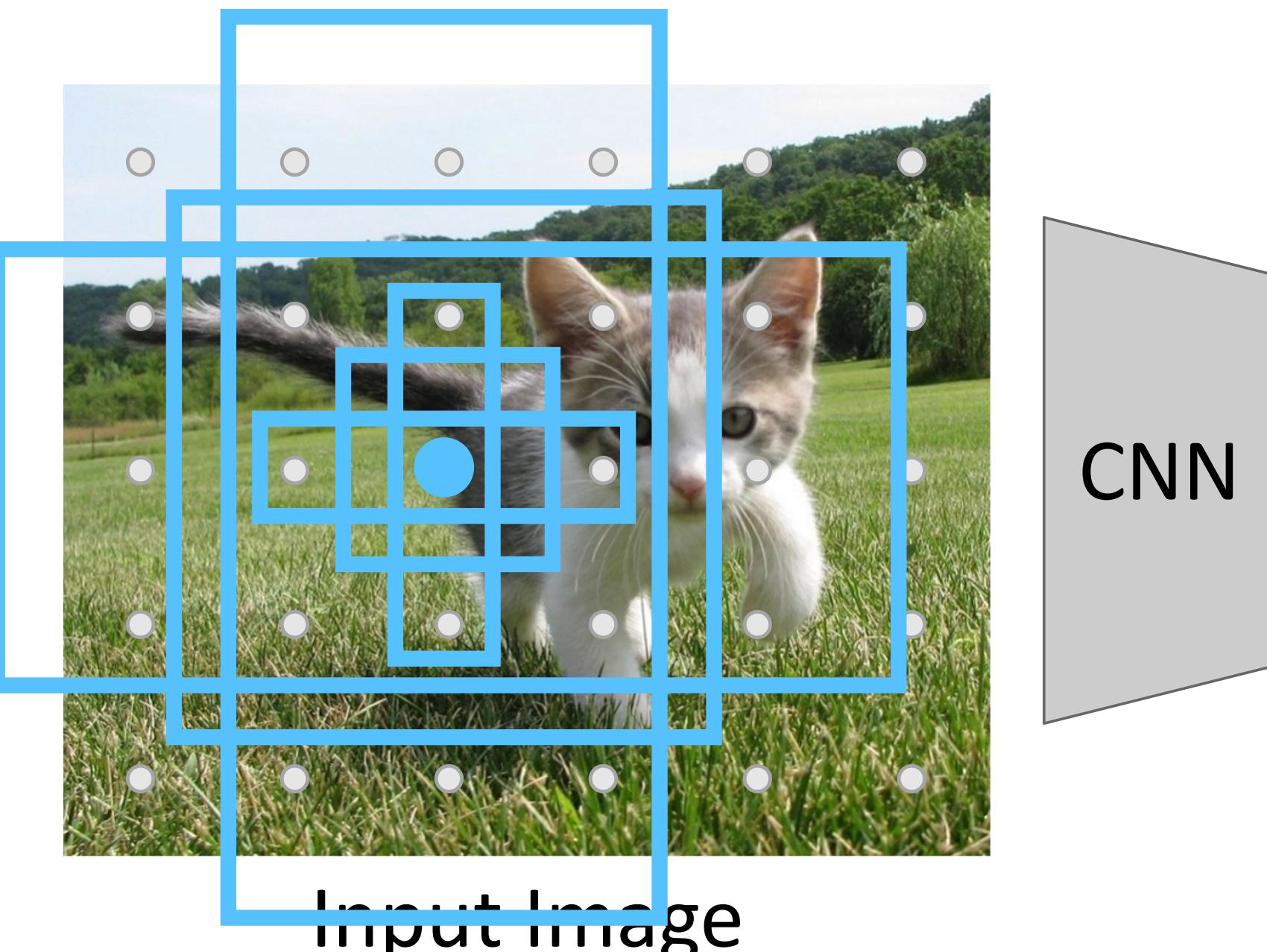
In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here $K = 6$)



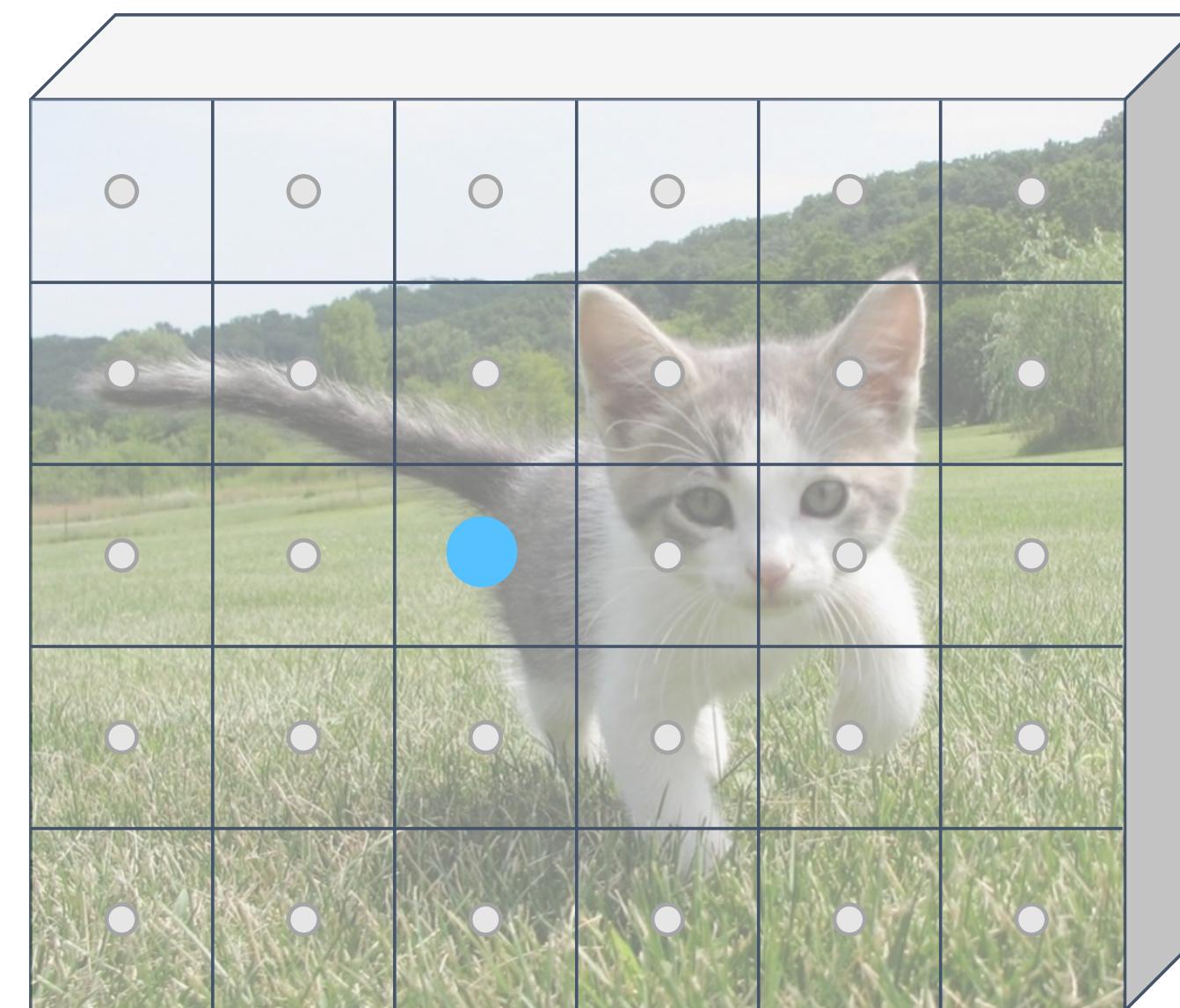
Neutral anchors: between 0.3 and 0.7 IoU with all GT boxes; ignored during training

Region Proposal Network (RPN)

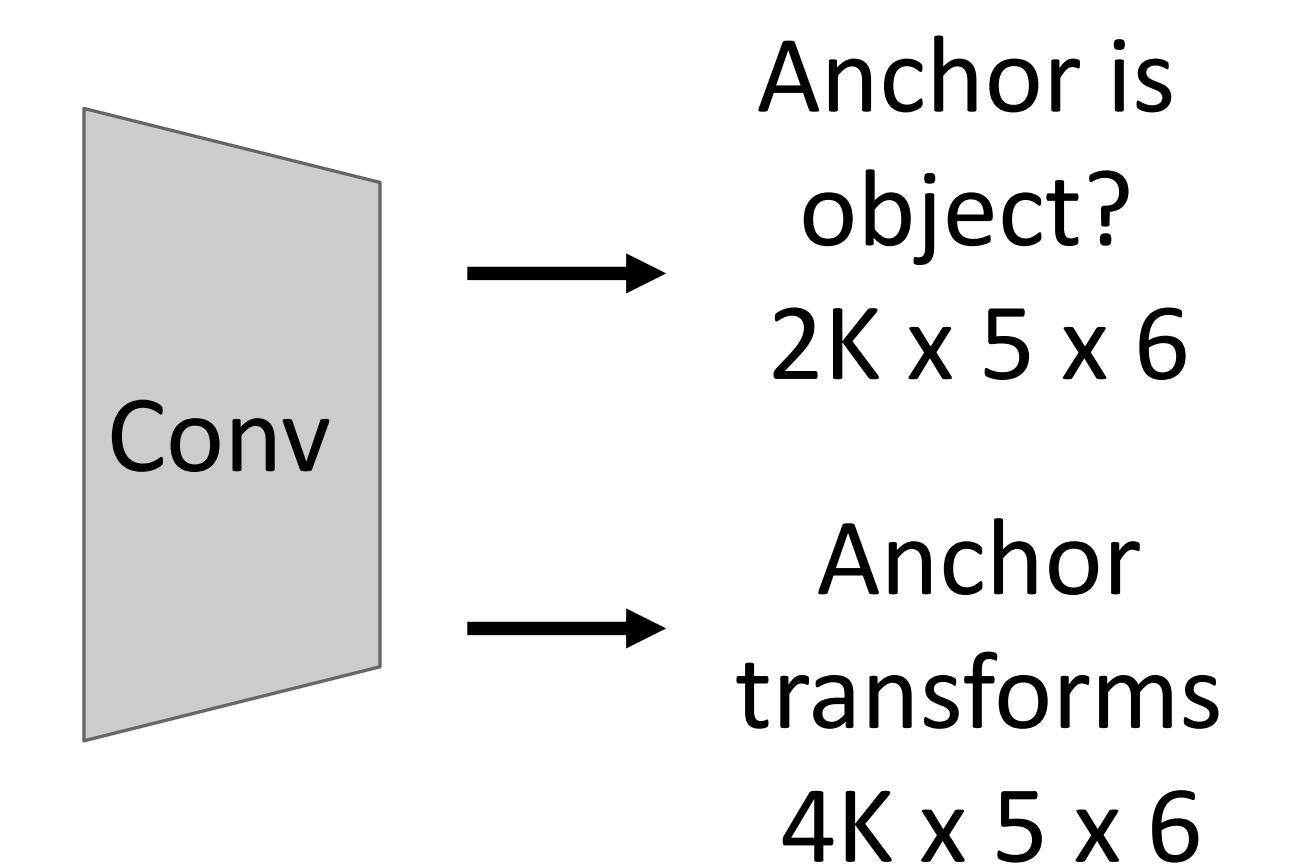
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here $K = 6$)

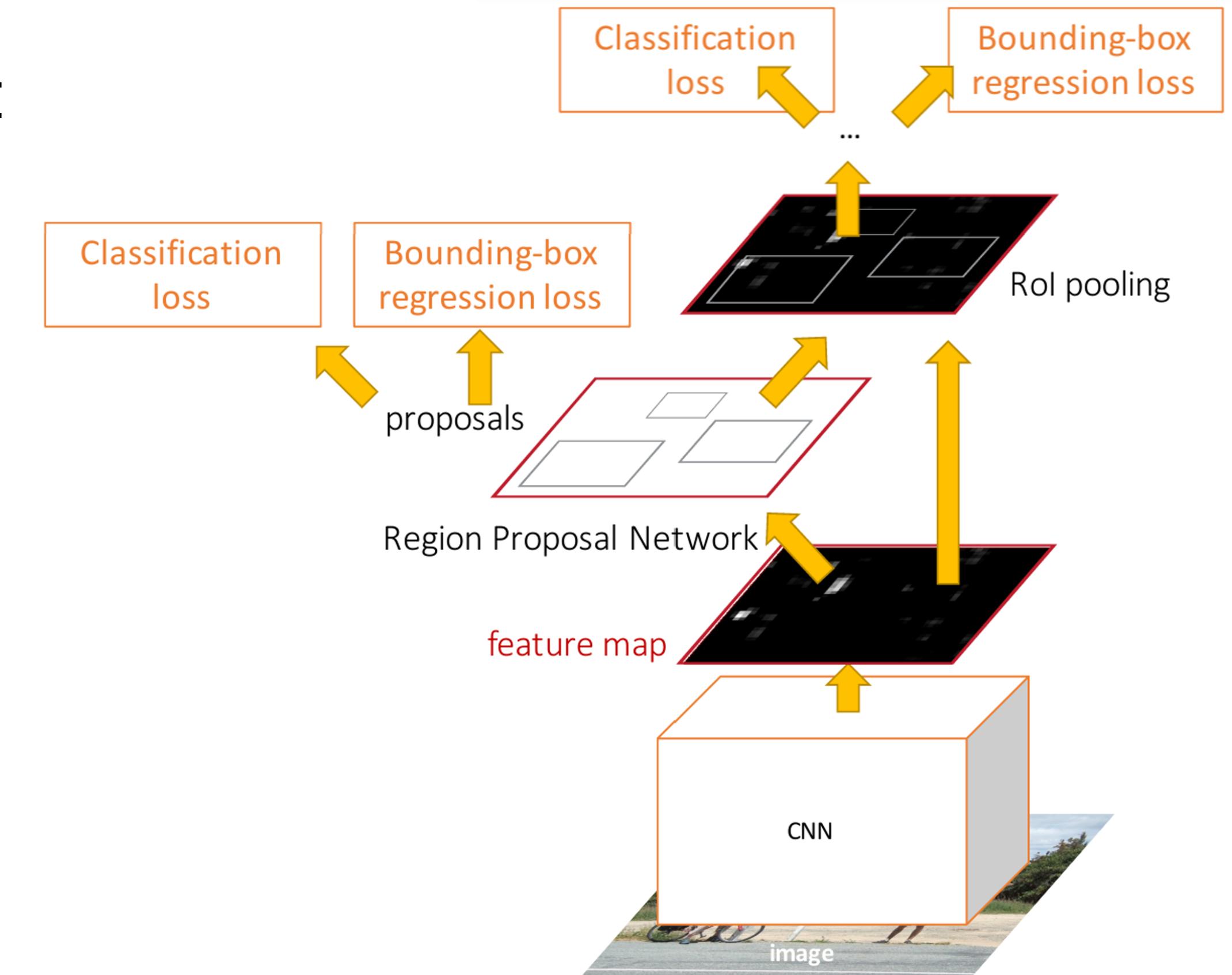


At test-time, sort all $K*5*6$ boxes by their positive score, take top 300 as our region proposals

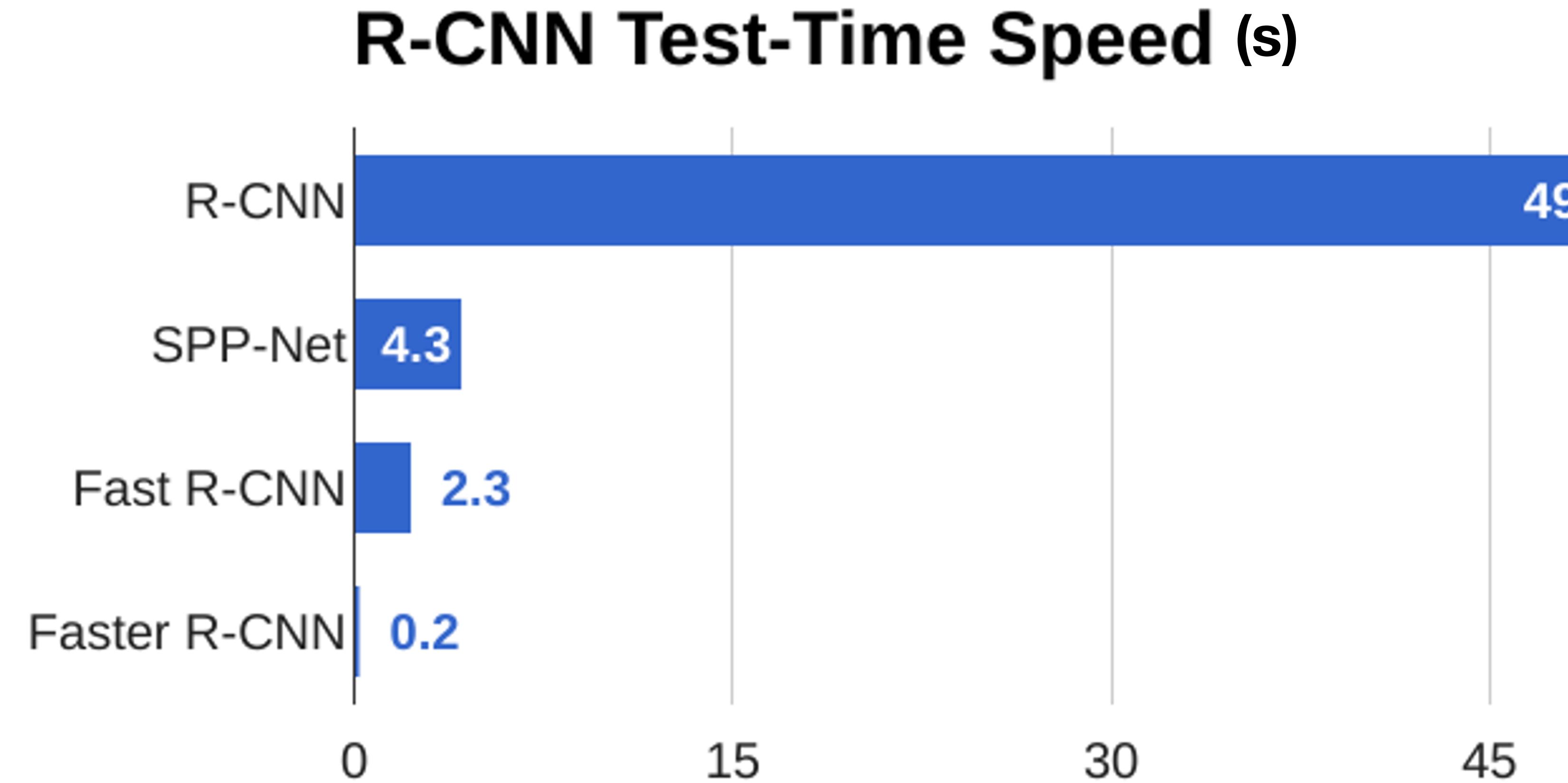
Faster R-CNN: Learnable Region Proposals

Jointly train four losses:

1. **RPN classification:** anchor box is object / not an object
2. **RPN regression:** predict transform from anchor box to proposal box
3. **Object classification:** classify proposals as background / object class
4. **Object regression:** predict transform from proposal box to object box



Faster R-CNN: Learnable Region Proposals



Extend Faster R-CNN to Image Segmentation: Mask R-CNN

Classification



“Chocolate Pretzels”

No spatial extent



Semantic Segmentation



Chocolate Pretzels,
Shelf

No objects, just pixels

Object Detection



Flipz, Hershey's, Keeze's

Multiple objects

Instance Segmentation



Extend Faster R-CNN to Instance Segmentation: Mask R-CNN

Instance Segmentation

Detect all objects in the image and identify the pixels that belong to each object (Only things!)

Approach

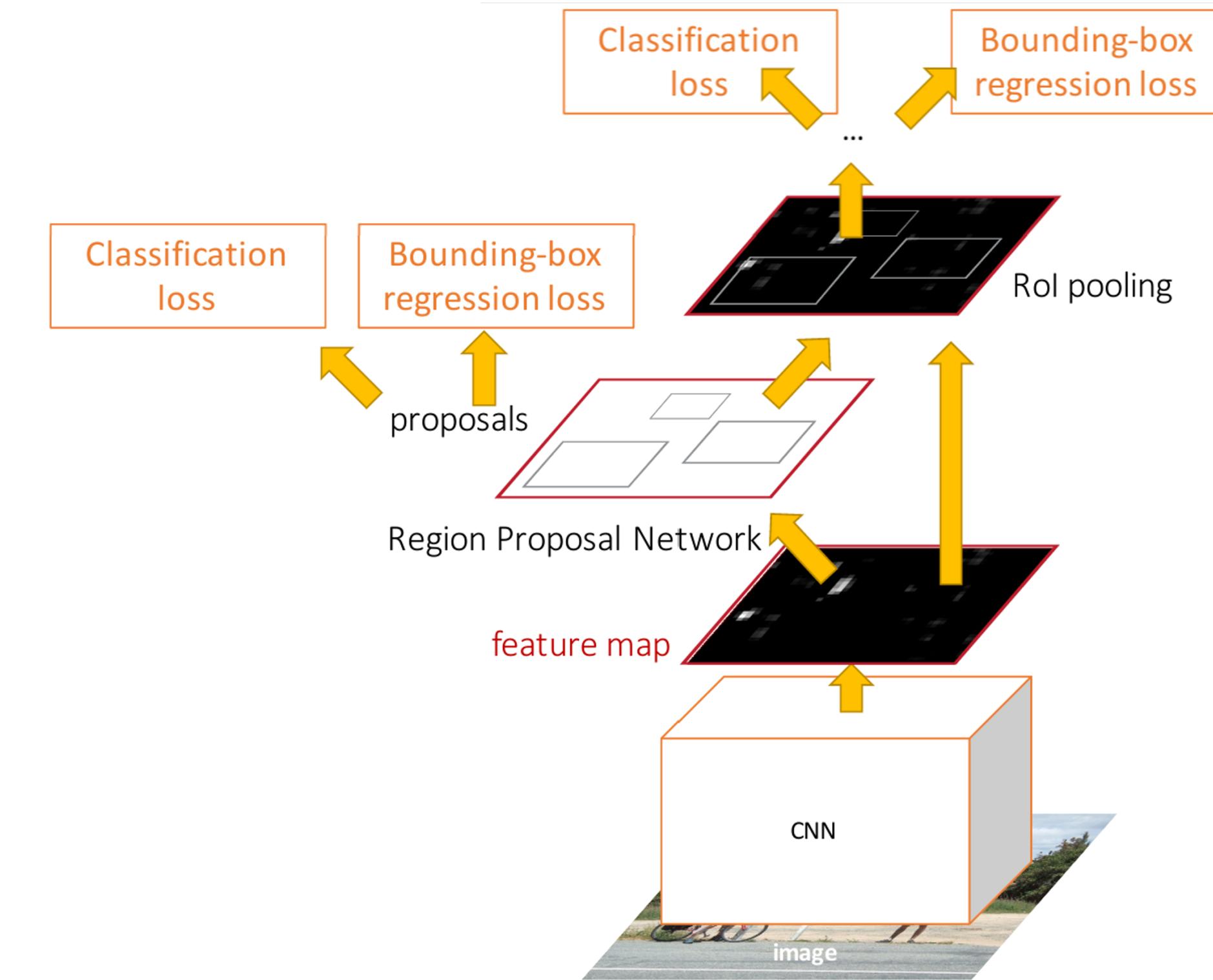
Perform object detection then predict a segmentation mask for each object detected!



Extend Faster R-CNN into Mask R-CNN

Faster R-CNN

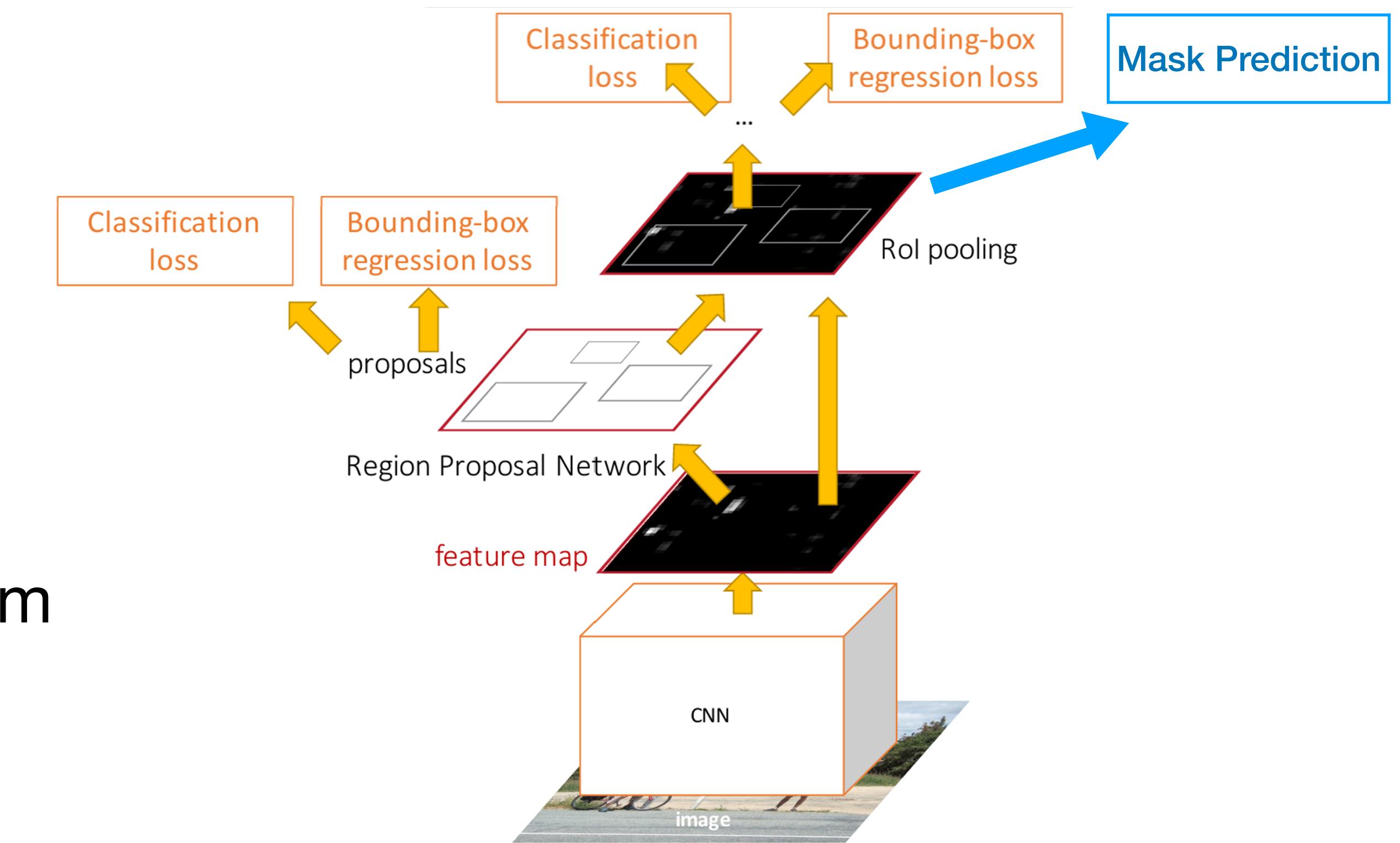
1. **Feature Extraction** at the image-level
2. **Regions of Interest** proposal from feature map
3. **In Parallel**
 1. **Object classification:** classify proposals
 2. **Object regression:** predict transform from proposal box to object box



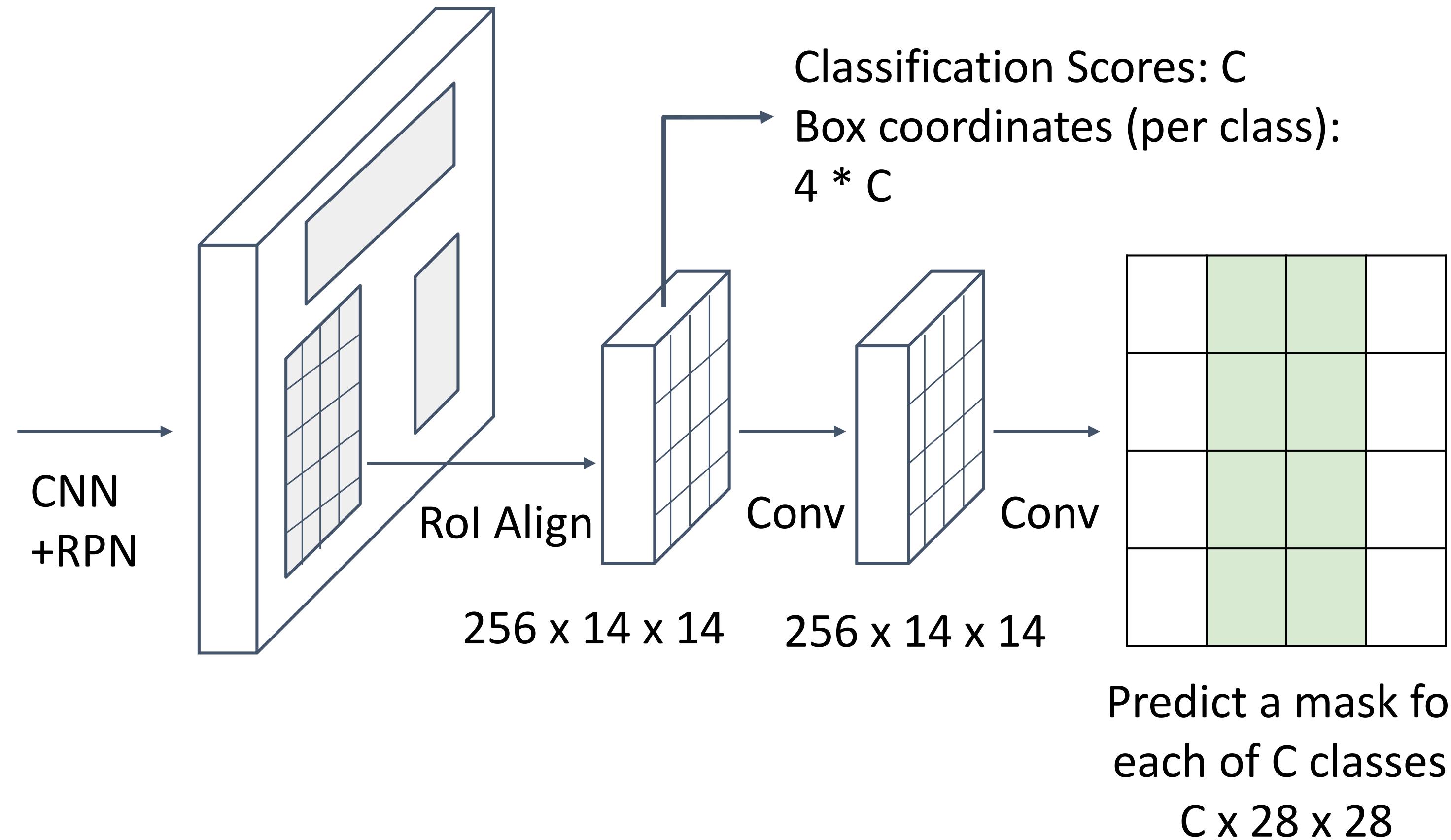
Extend Faster R-CNN into Mask R-CNN

Mask R-CNN

1. **Feature Extraction** at the image-level
2. **Regions of Interest** proposal from feature map
3. **In Parallel**
 - a. **Object Classification:** classify proposals
 - b. **Object Regression:** predict transform from proposal box to object box
 - c. **Mask Prediction:** predict a binary mask for every region



Mask R-CNN



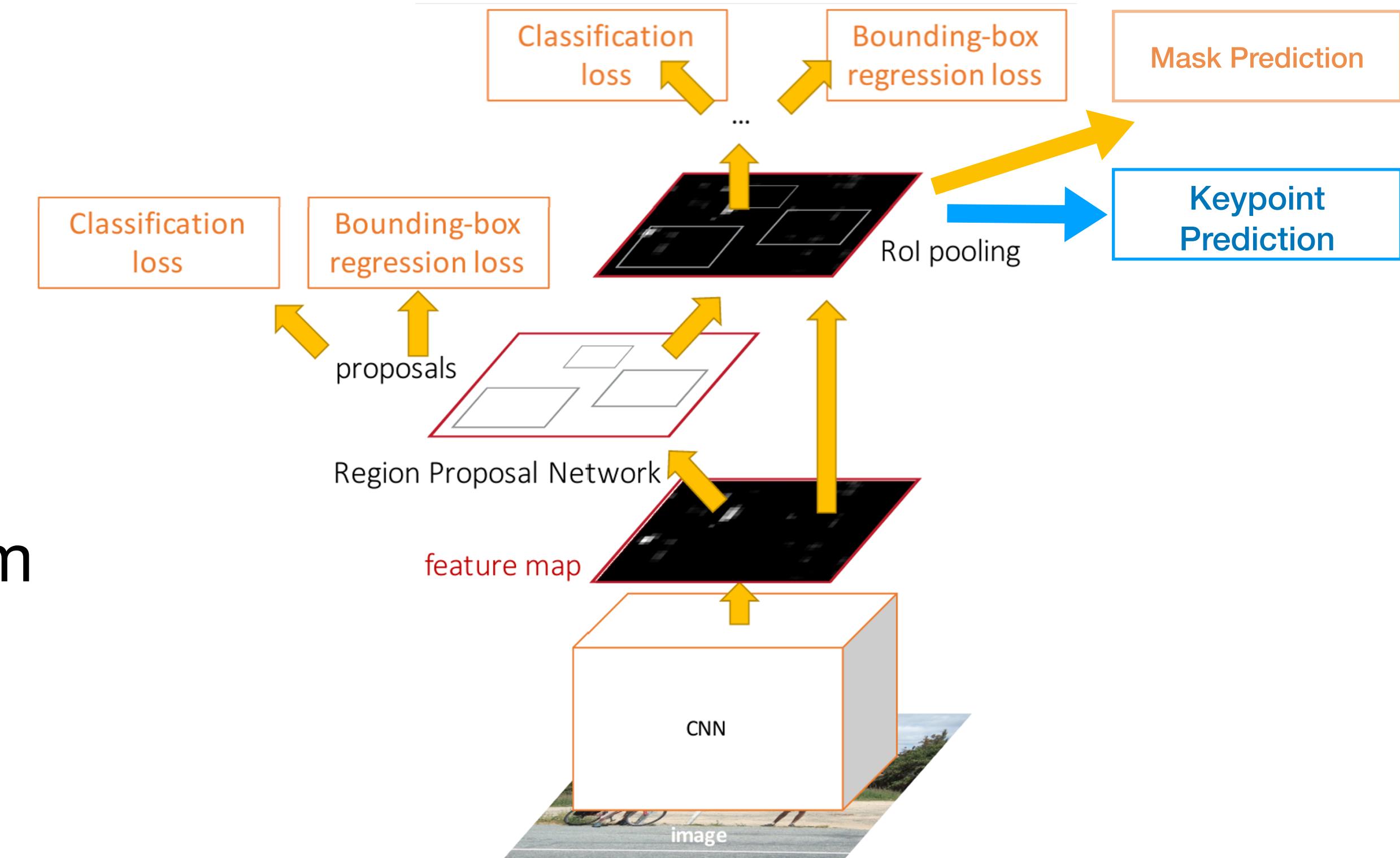
Mask R-CNN: Very Good Results!



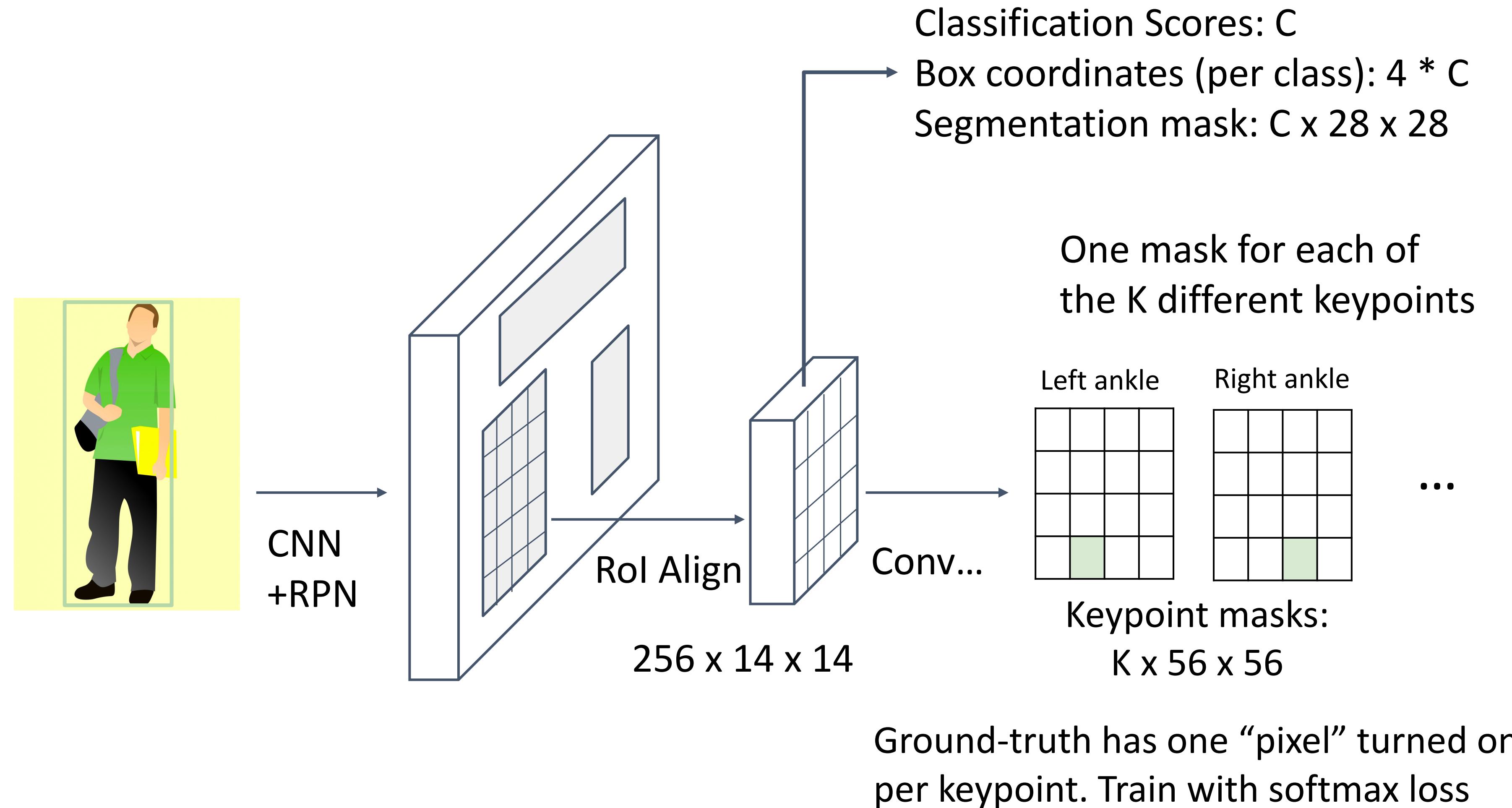
Mask R-CNN for Human Pose Estimation

Mask R-CNN

1. **Feature Extraction** at the image-level
2. **Regions of Interest** proposal from feature map
3. **In Parallel**
 - a. **Object Classification:** classify proposals
 - b. **Object Regression:** predict transform from proposal box to object box
 - c. **Mask Prediction:** predict a binary mask for every region
 - d. **Keypoint Prediction:** predict binary mask for human key points



Mask R-CNN for Human Pose Estimation



Mask R-CNN for Human Pose Estimation



Two Stage vs One Stage Detectors

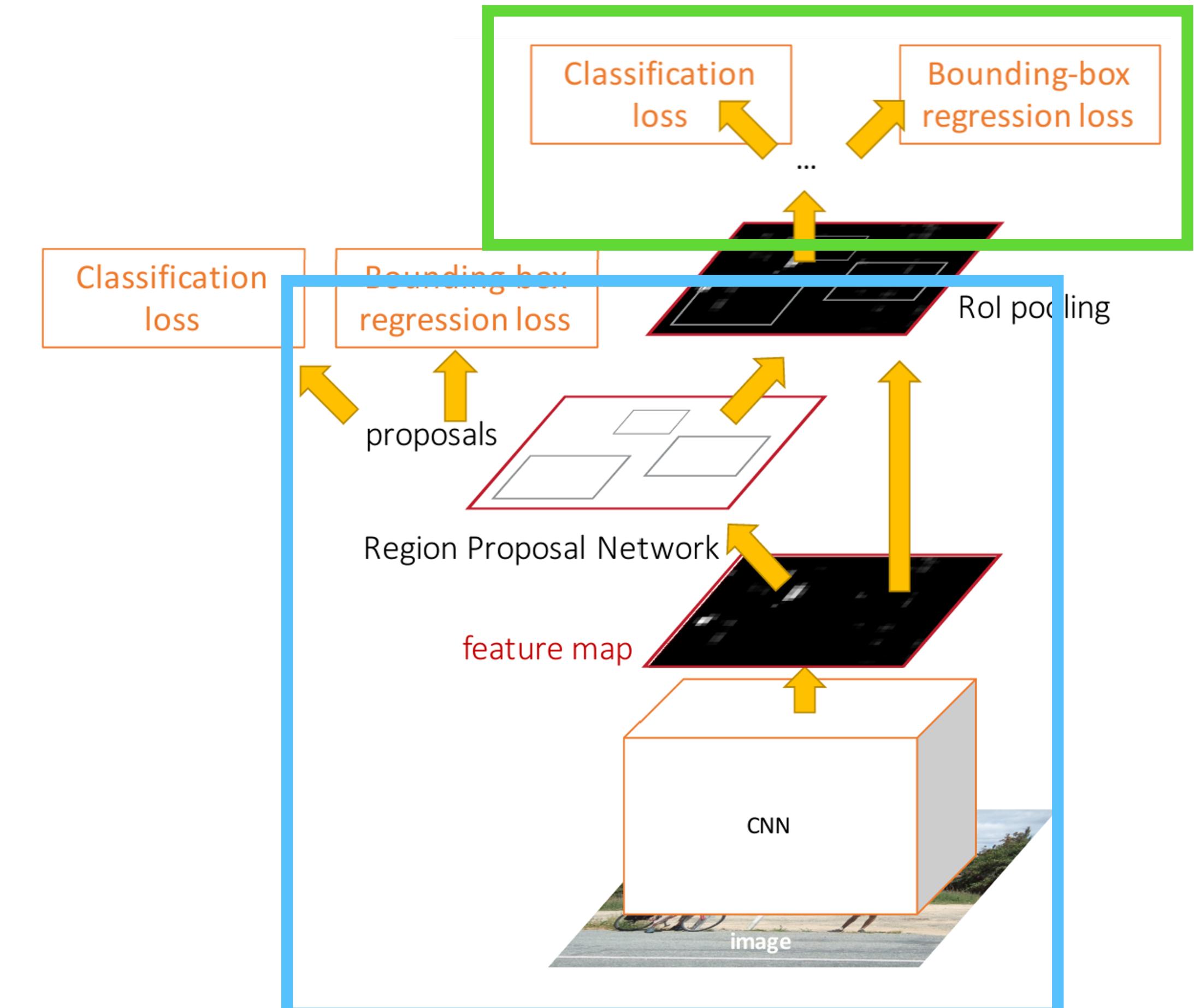
Faster R-CNN is a **two-stage object detector**

First stage: Run once per image

- Backbone Network
- Region Proposal Network

Second stage: Run once per region

- Crop features: RoI pool / align
- Predict Object Class
- Prediction bbox offset



DR

DeepRob

Lecture 13
Object Detectors and Segmentation
University of Michigan and University of Minnesota

